



UNIVERSITÉ  
DE GENÈVE

FACULTY OF SCIENCE  
Department of Informatics



stochastic  
information  
processing

□

# Information Theory for Data Science and Machine Learning

Slava Voloshynovskiy

# Miscellanea

---

## Lecturer:

Prof. **Slava Voloshynovskiy**, Stochastic Information Processing (SIP)  
Group, Department of Computer Science, UNIGE  
**E-mail:** [svolos@unige.ch](mailto:svolos@unige.ch); **Office:** 215

## Teaching Assistants:

**Guillaume Quetant**, PhD student, Stochastic Information Processing (SIP)  
Group, Department of Computer Science, UNIGE  
**E-mail:** Guillaume.Quetant@unige.ch

**Mariia Drozdova**, PhD student, Stochastic Information Processing (SIP)  
Group, Department of Computer Science, UNIGE  
**E-mail:** mariia.drozdova@unige.ch

# Organization and evaluation

---

## Lectures

- Friday 10:00-12:00, Bat A/404-407, starting from February 25, 2022

## HWs

Wednesday 8:00-10:00, Bat D/Amphi, starting from March 3, 2023

**Web Site:** [moodle 12x004](#) Théorie de l'information pour la science des données et l'apprentissage automatique

## Organization

- One semester
  - 14 lectures
  - 10 HWs

## The final note is composed of

- Oral exam or CC (consisting of two parts during the semester) - 2/3
- HWs - 1/3

# Copyright notice

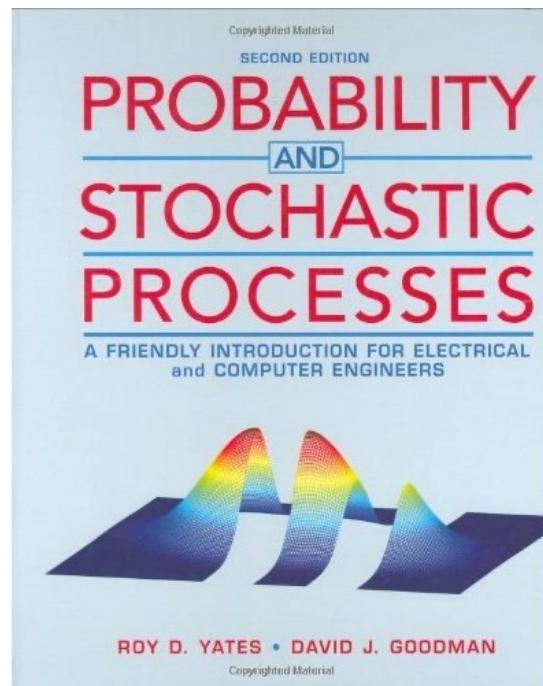
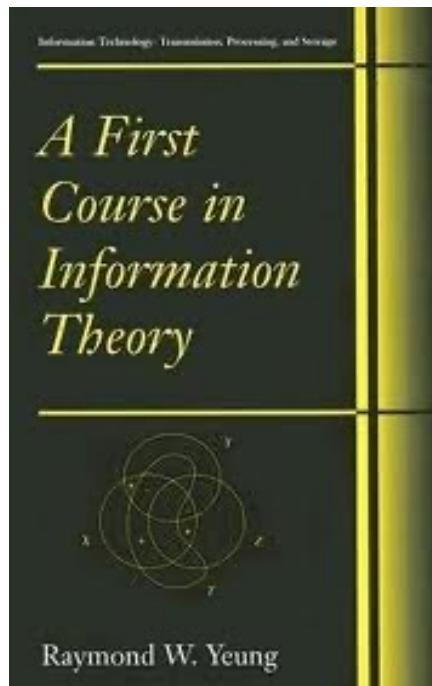
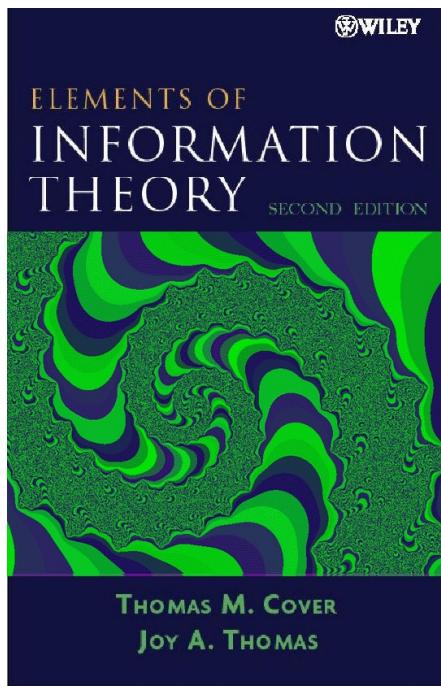
---

This course uses in part some materials from previously mentioned recommended sources.

The usage of the slides in commercial or educational purposes  
is prohibited without authorization of the course leaders and  
permission of the above document authors and copyright owners.

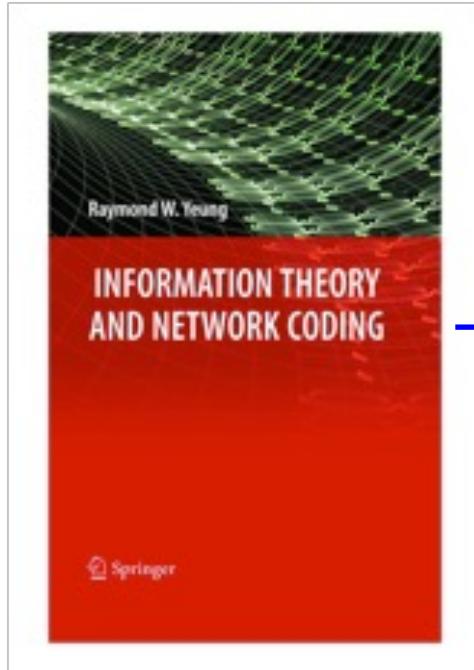
# Recommended books

---



# Recommended books

---



Free download as

Raymond W. Yeung

Information Theory and  
Network Coding

SPIN Springer's internal project number, if known

January 28, 2008

<http://iest2.ie.cuhk.edu.hk/~whyeung/book2/>

<http://iest2.ie.cuhk.edu.hk/~whyeung/post/manuscript/main2.pdf>

# Content of course

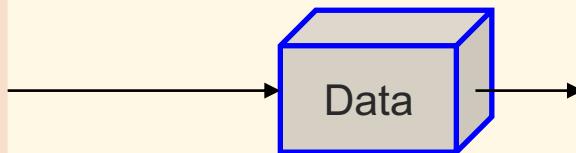
---

- Theme 1 – Basic statistical data models
- Theme 2 – Information theoretic measures

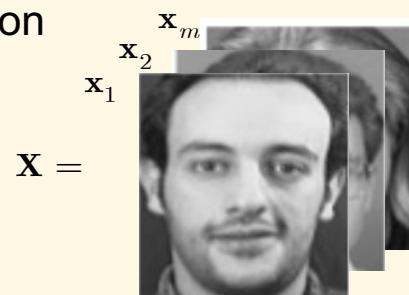
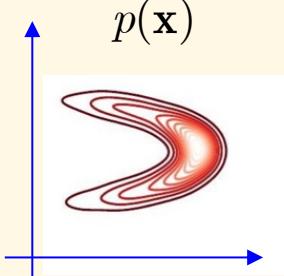
# Scope of class

## IT measures

Statistical source of data



Complex statistical distribution



## Data Science and Machine Learning

- Data analysis/visualization
- Data processing
  - clustering
  - denoising
  - restoration
  - classification
- Data compression
- Data transmission
- Data generation

# Content of this lecture

---

**In this lecture we will:**

- Introduce basic notations
- Recall basics of probability and matrix algebra
- Consider basic data models
- Recall transformation of random variables and vectors
- Recall weak law of large numbers and the central limit theorem (CLT)

# Basic notions

---

- Random variables and vectors
  - Pmf/pdf
  - Moments
- Data representation
- Modeling of random vectors
  - Chain rule for probability
  - Independence and histograms
  - Markov chains
  - Conditional independence
  - Marginalization

# Notations

---

- **Basic notations of scalars and vectors**

|        | Realization        | Random             |
|--------|--------------------|--------------------|
| Scalar | $x$                | $X$                |
| Vector | $\mathbf{x} = x^N$ | $\mathbf{X} = X^N$ |

$$x_i \in \mathcal{X} \quad \text{or} \quad x_i \in S_X$$

- Discrete alphabets

$$\mathcal{X} = \mathbb{Z} \quad \text{- integer}$$

$$\mathcal{X} = \{0,1\} \quad \text{- binary}$$

$$\mathcal{X} = \{a,b,\dots,z\}$$

$$\mathcal{X} = \{January, February, \dots, December\}$$

- Continuous alphabets

$$\mathcal{X} = \mathbb{R}$$

- Vector  $\mathbf{x} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{array}{c} 1 \\ 2 \\ \vdots \\ N \end{array} \begin{array}{c} \textcircled{1} \\ \textcircled{2} \\ \vdots \\ \textcircled{N} \end{array}$

$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{X} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

# Notations

---

## Matrix operations

|                     |                                    |
|---------------------|------------------------------------|
| $\mathbf{x}^T$      | Transposed vector $\mathbf{x}$     |
| $\mathbf{I}$        | Identity matrix                    |
| $\mathbf{A}^T$      | Transposed matrix $\mathbf{A}$     |
| $\mathbf{A}^{-1}$   | Inverse matrix $\mathbf{A}$        |
| $tr(\mathbf{A})$    | Trace of matrix $\mathbf{A}$       |
| $\det  \mathbf{A} $ | Determinant of matrix $\mathbf{A}$ |

Recall B: matrix operations

# Random variables

---

- **Basic notations of probability**

$\Pr[X = x]$  Probability of event

- Discrete random variables

$X \sim p_X(x)$  Probability mass function (**pmf**) of discrete random variable  $X$   
(or  $X \sim p(x)$ )

$p_X(x | Y = y)$  Conditional pmf of discrete random variable  $X$  given  $Y$   
(or  $p_X(x | y)$ )

- Continuous random variables

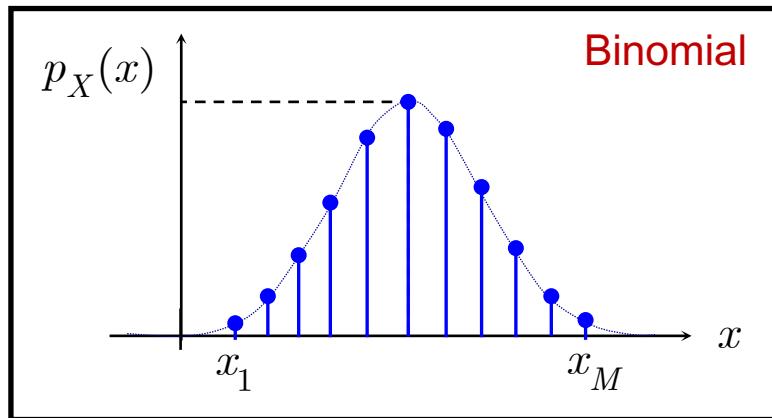
$X \sim f_X(x)$  Probability density function (**pdf**) of continuous random variable  $X$   
(or  $X \sim f(x)$ )

$f_X(x | Y = y)$  Conditional pdf of continuous random variable  $X$  given  $Y$   
(or  $f_X(x | y)$ )

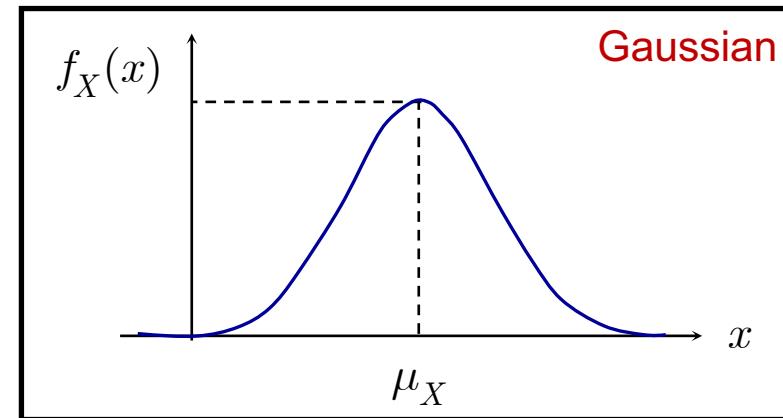
# Random variables: pmf and pdf

- Basic notations of probability

- Discrete random variables



- Continuous random variables



$$p_X(x) = \Pr[X = k] = \binom{M}{k} p^k (1 - p)^{M-k}$$
$$k = 0, 1, \dots, M$$
$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma_X^2}} e^{-\frac{(x-\mu_X)^2}{2\sigma_X^2}}$$

$$X \sim p_X(x) = \mathcal{B}(M, p)$$

$$X \sim f_X(x) = \mathcal{N}(\mu_X, \sigma_X^2)$$

# Random variables: moments

## ▪ Moments

### ▪ Discrete random variables

1<sup>st</sup> moment  
(mean)

$$\mu_X = E_{p_X(x)}[X] = \sum_{i=1}^M xp(x_i)$$

### ▪ Continuous random variables

$$\mu_X = E_{f_X(x)}[X] = \int_{x \in \mathcal{X}} xf_X(x) dx$$

2<sup>nd</sup> moment  
(variance)

$$\sigma_X^2 = Var[X] = E_{p_X(x)}[|X - \mu_X|^2]$$

$$\sigma_X^2 = Var[X] = E_{f_X(x)}[|X - \mu_X|^2]$$

$$= \sum_{i=1}^M (x_i - \mu_X)^2 p(x_i)$$

$$= \int_{x \in \mathcal{X}} (x - \mu_X)^2 f_X(x) dx$$

3<sup>rd</sup> moment  
(skewness)

$$E_{p_X(x)}[|X - \mu_X|^3] / \sigma_X^3$$

$$E_{f_X(x)}[|X - \mu_X|^3] / \sigma_X^3$$

4<sup>th</sup> moment  
(kurtosis)

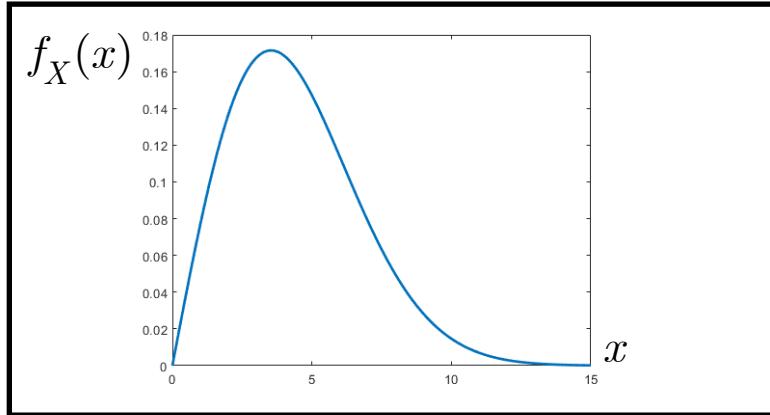
$$E_{p_X(x)}[|X - \mu_X|^4] / \sigma_X^4$$

$$E_{f_X(x)}[|X - \mu_X|^4] / \sigma_X^4$$

# Random variables: moments

- General case: given a pmf/pdf, we can compute moments (not for all)

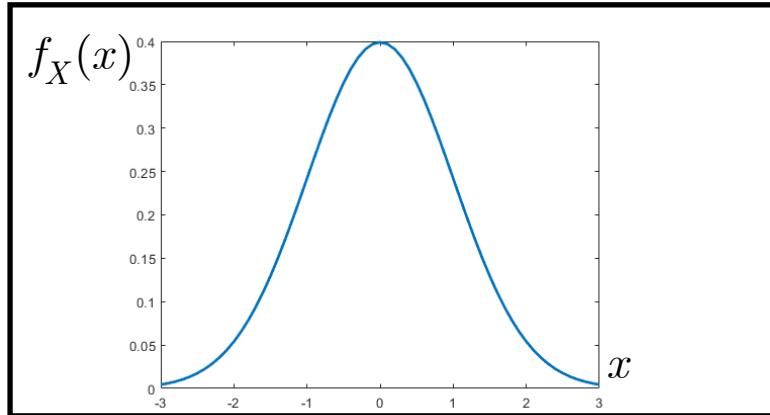
Any pdf



$$\begin{bmatrix} E_{p_X(x)}[X] \\ E_{p_X(x)}[|X - \mu_X|^2] \\ E_{p_X(x)}[|X - \mu_X|^4] \\ E_{p_X(x)}[|X - \mu_X|^5] \\ \vdots \end{bmatrix}$$

A red X is placed over the third row of the matrix.

$\mathcal{N}(\mu_X, \sigma_X^2)$



$$\begin{bmatrix} E_{p_X(x)}[X] \\ E_{p_X(x)}[|X - \mu_X|^2] \end{bmatrix}$$

Intuitive interpretation

# Basic operation: moments of sum of two variables

---

- Expected value of function of two random variables  $Z = g(X, Y)$

$$E[Z] = \sum_{x \in S_X} \sum_{y \in S_Y} g(x, y) p_{X,Y}(x, y)$$

$$Z = g(X, Y) = X + Y$$

- **Mean** of sum of two random variables

$$E[Z] = E[X + Y] = E[X] + E[Y]$$

- **Variance** of sum of two random variables

$$Var[Z] = Var[X + Y] = Var[X] + Var[Y] + 2E[(X - \mu_X)(Y - \mu_Y)]$$

Remark: we will see pdf of sum of two and many variables later

# Recall of variance and its properties

---

- Proof

$$\begin{aligned}Var[Z] &= Var[X + Y] = E\left[\left((X + Y) - (\mu_X + \mu_Y)\right)^2\right] \\&= E\left[\left((X - \mu_X) + (Y - \mu_Y)\right)^2\right] \\&= E\left[\left((X - \mu_X)^2 + 2(X - \mu_X)(Y - \mu_Y) + (Y - \mu_Y)^2\right)\right] \\&= E\left[(X - \mu_X)^2\right] + 2E\left[(X - \mu_X)(Y - \mu_Y)\right] + E\left[(Y - \mu_Y)^2\right] = \\&= Var[X] + Var[Y] + 2E\left[(X - \mu_X)(Y - \mu_Y)\right]\end{aligned}$$

—

# Random variables: covariance

---

- Covariance between two random variables

$$\begin{aligned} \text{Cov}[X, Y] &= E_{p(x,y)}[(X - \mu_X)(Y - \mu_Y)] = \sum_{(x,y) \in S_{X,Y}} p(x,y)(x - \mu_X)(y - \mu_Y) \\ &= E_{p(x,y)}[XY] - \mu_X\mu_Y = \sigma_{XY} \end{aligned}$$

- Variance of sum via covariance

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$

- **Uncorrelated** random variables

$$E[XY] = E[X]E[Y] = \mu_X\mu_Y$$

$$\text{Cov}[X, Y] = \mu_X\mu_Y - \mu_X\mu_Y = 0$$

$$\Rightarrow \text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$$

# Random variables: covariance

---

- Correlation coefficient

$$\rho_{XY} = \frac{Cov[X, Y]}{\sqrt{Var[X]Var[Y]}} = \frac{\sigma_{XY}}{\sigma_X\sigma_Y}$$

$$-1 \leq \rho_{XY} \leq 1$$

# Random variables: summary

---

- Properties of variance and covariance
  - Let  $X$  and  $Y$  be random variables and  $a, b, c$  be constants.

$$Var[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$$0 \leq Var[X] \leq E[X^2]$$

$$Var[cX] = c^2Var[X]$$

$$Var[aX + bY] = a^2Var[X] + 2abCov[X, Y] + b^2Var[Y]$$

$$Cov[X, Y] = E_{p(x,y)}[(X - \mu_X)(Y - \mu_Y)] = \sum_{(x,y) \in S_{X,Y}} p(x,y)(x - \mu_X)(y - \mu_Y)$$

$$= E_{p(x,y)}[XY] - \mu_X\mu_Y = \sigma_{XY}$$

See **Appendix A** for the recall, if needed

- From two variables to multidimensional vectors

# Basic notions

---

- Random variables and vectors
  - Pmf/pdf
  - Moments
- Data representation
- Modeling of random vectors
  - Chain rule for probability
  - Independence and histograms
  - Markov chains
  - Conditional independence
  - Marginalization

# Random vectors

---

- Vector  $\mathbf{x} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \xrightarrow{\text{random}} \mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{X} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$

$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{X} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_i \\ \vdots \\ X_N \end{bmatrix}$$

$p_X(x_i)$   
 $f_X(x_i)$

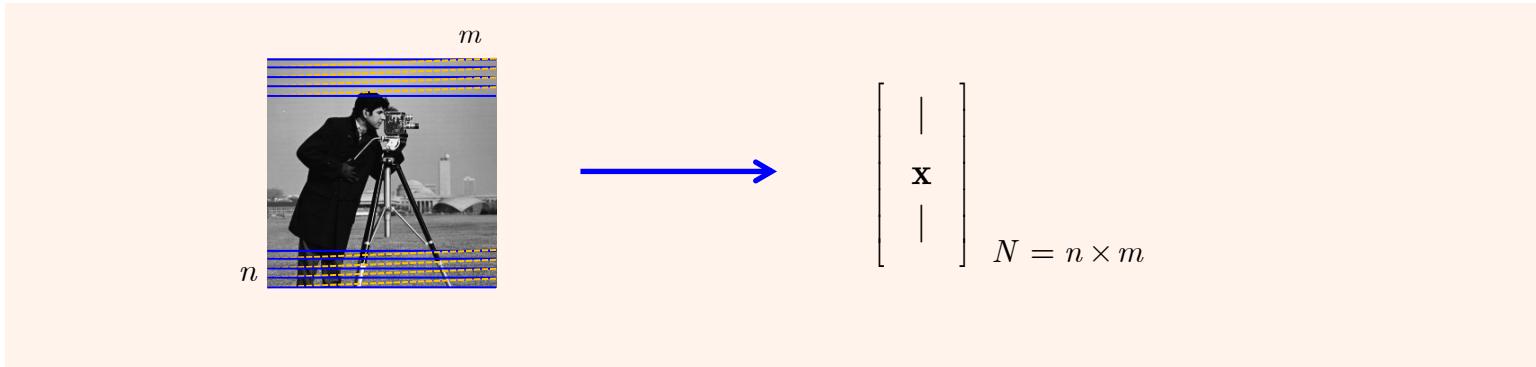
- Joint pmf or pdf

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{x}}(x_1, x_2, \dots, x_N)$$

$$f_{\mathbf{x}}(\mathbf{x}) = f_{\mathbf{x}}(x_1, x_2, \dots, x_N)$$

# Data representation - 1

- **Example:**
  - Consider an image as a matrix
  - Read in a zig-zag order into one vector



- Joint pmf defines **all statistical relations** between all elements of image/vector

$$\mathbf{X} = \begin{bmatrix} | & \\ \mathbf{x} & | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

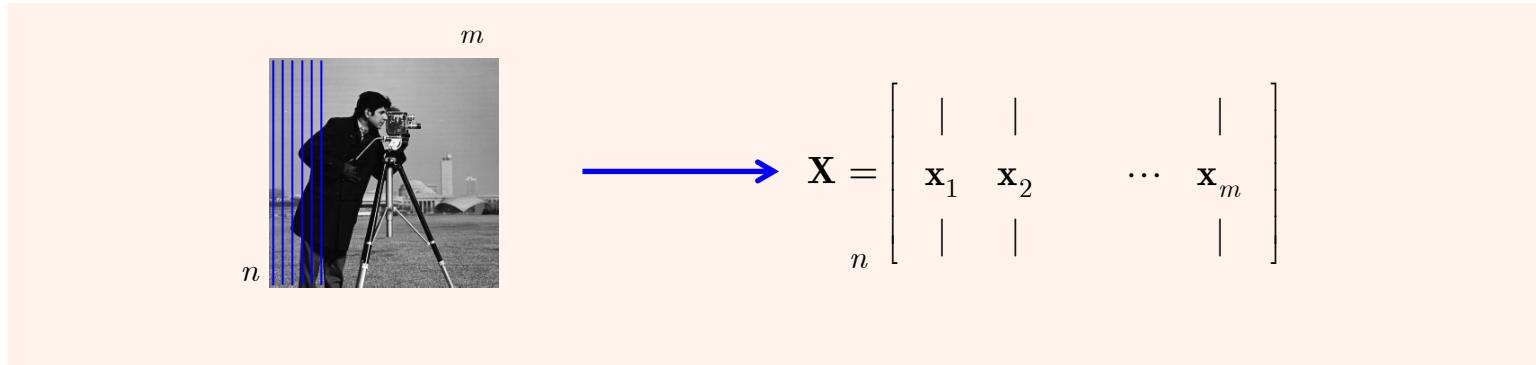
Diagram illustrating the joint pmf of an image vector  $\mathbf{X}$ . The vector is shown as a column of elements  $X_1, X_2, \dots, X_N$ . Blue arrows indicate the mapping from the image pixels to the corresponding elements in the vector.

- Complexity
  - Modern image of size  $1000 \times 1000$   
 $N = n \times m = 1000 \times 1000 = 1000000$   
All pair-wise, triple-wise, etc.



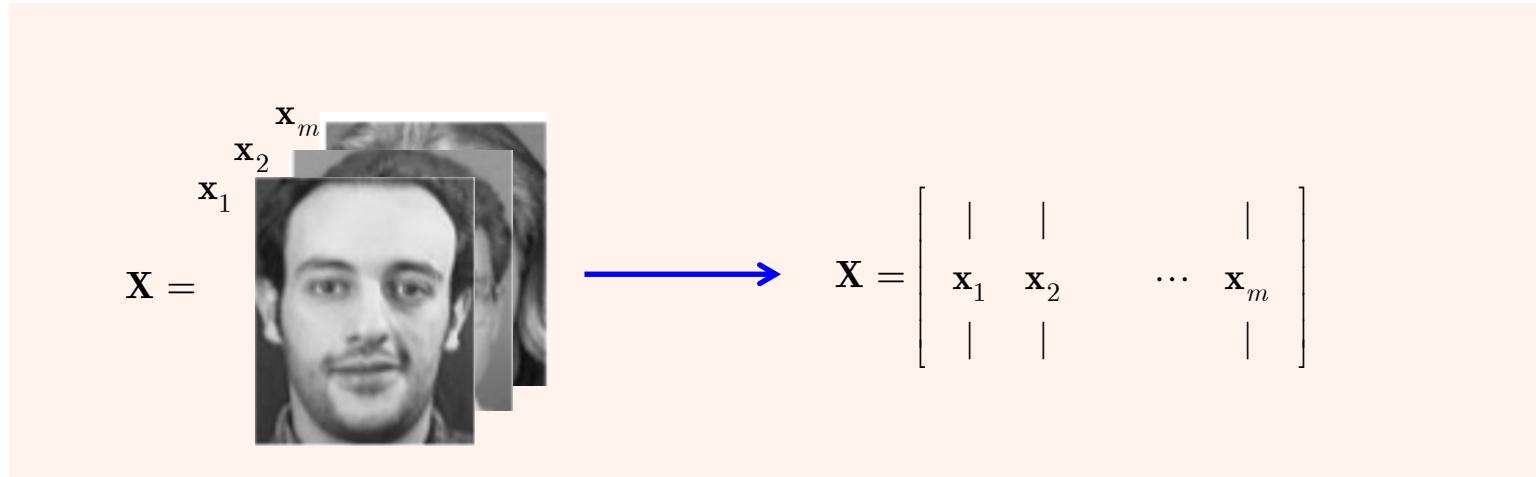
# Data representation - 2

- **Example:**
  - Consider an image as a matrix
  - Read by rows as a matrix



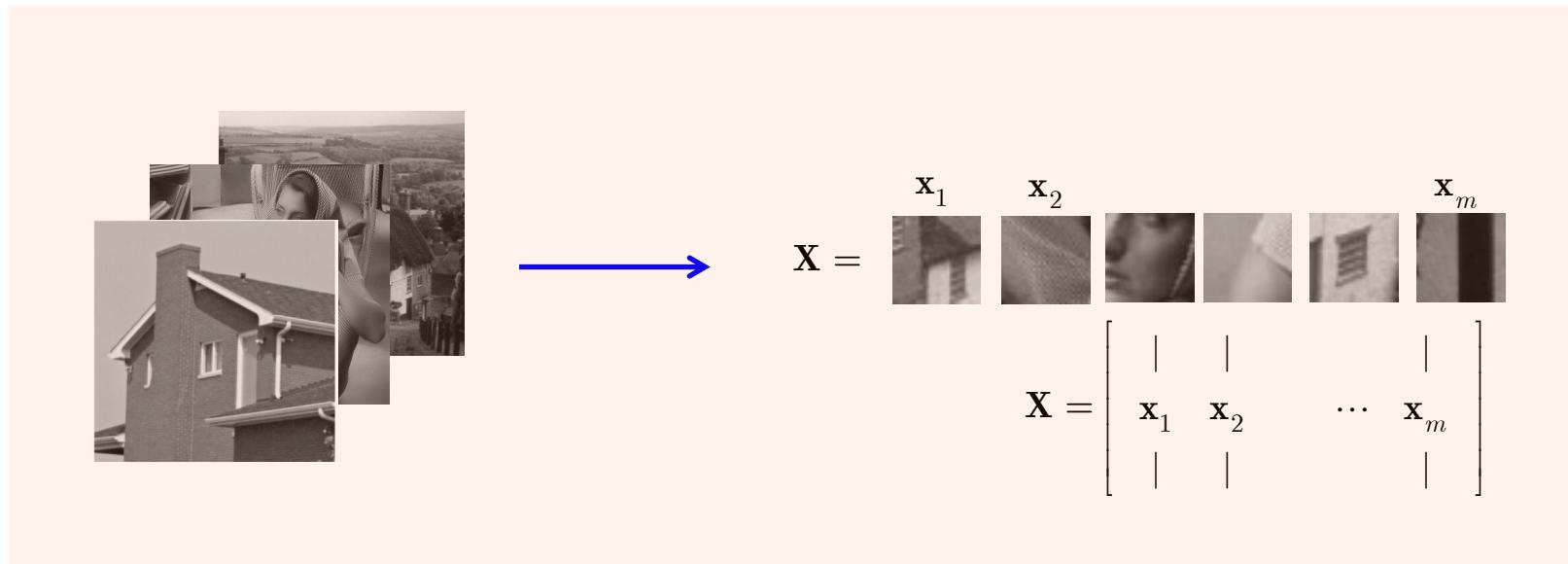
# Data representation - 3

- Example:
  - image represents a column in a data matrix



# Data representation - 4

- Example
  - image representation as patches



# Basic notions

---

- Random variables and vectors
  - Pmf/pdf
  - Moments
- Data representation
- Modeling of random vectors
  - Chain rule for probability
  - Independence and histograms
  - Markov chains
  - Conditional independence
  - Marginalization

# Random vectors: models

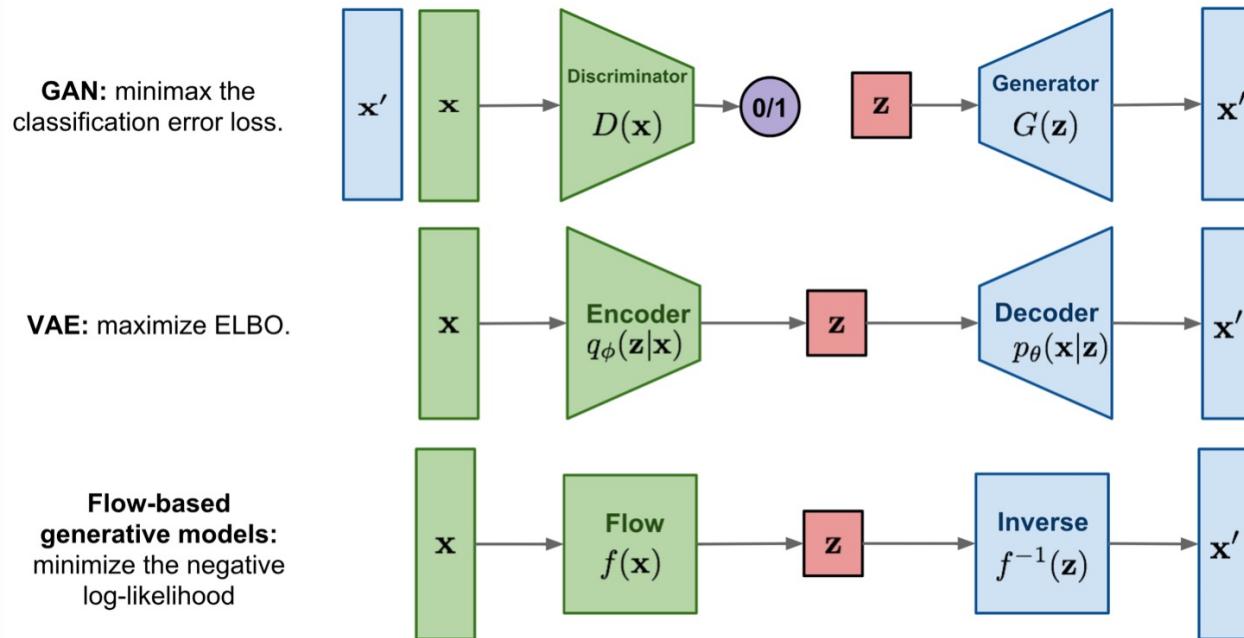
---

- Modeling of real data represented by random vectors is a very complex problem
  - Computational issues
  - Memory storage
  - Variability between data in various applications (medicine, astro, biometrics, etc.)
- Existing approaches to attack this issue are based on different assumptions:
  - Models based on independence
  - Models based on low-order dependencies /Markov chain/
  - Models based on second-order statistics /multivariate Gaussians/
  - Advanced models based on multi-level, multi-stage representations /deep networks/:
    - Variational autoencoder (VAE)
    - Generative adversarial networks (GAN)
    - Flow based models
    - Diffusion and score based models

# Random vectors: sota models

---

- State of the art (sota) models



However, to understand these complex models we need first to recall simple ones.

<https://lilianweng.github.io/lil-log/2018/10/13/flow-based-deep-generative-models.html>

# Random vectors: via transformation

- State of the art models

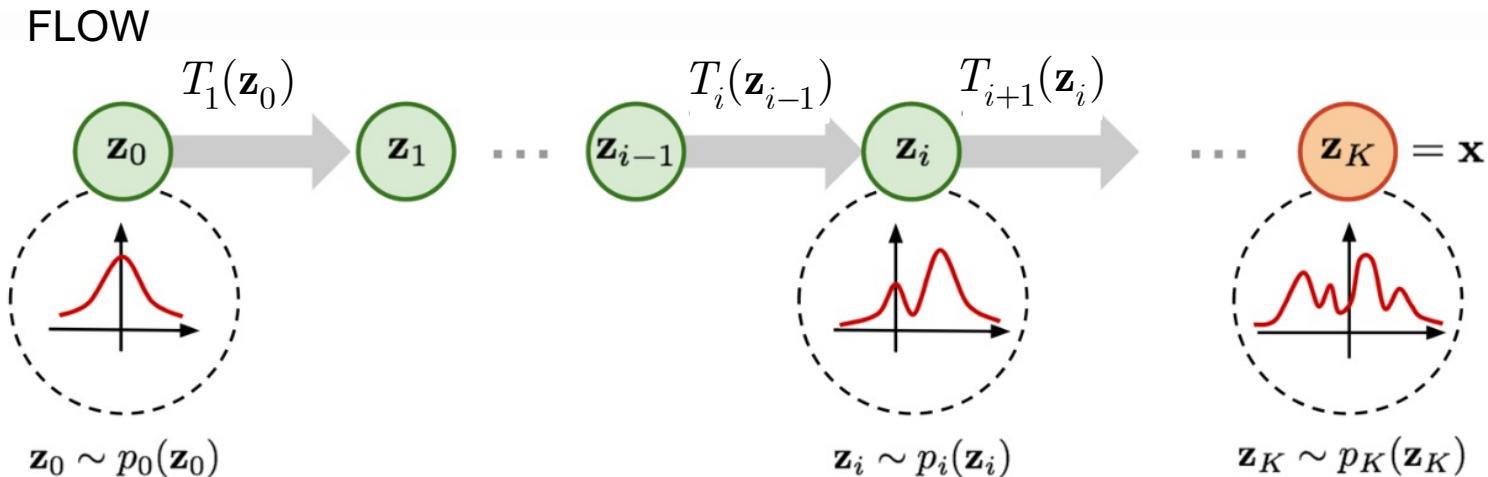


Illustration of a normalizing flow model, transforming a simple distribution to a complex one step by step

<https://lilianweng.github.io/lil-log/2018/10/13/flow-based-deep-generative-models.html>

# Random vectors: chain rule decomposition

- Chain rule for probability: joint probability

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{x}}(x_1, x_2, \dots, x_N) = p(x_1, x_2, \dots, x_N) \quad (\text{Notations})$$

$$\begin{aligned} p(x_1, x_2, \dots, x_N) &= p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots p(x_N|x_{N-1}, \dots, x_2, x_1) = \\ &= \prod_{i=1}^N p(x_i|x_{i-1}, \dots, x_2, x_1) \end{aligned}$$

$$p(x_1, x_2) = p(x_1)p(x_2|x_1) = p(x_2)p(x_1|x_2)$$

$$\Rightarrow p(x_2|x_1) = \frac{p(x_2)p(x_1|x_2)}{p(x_1)}$$

# Random vectors: independence

- Using the chain rule

$$\begin{aligned} p(x_1, x_2, \dots, x_N) &= p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots p(x_N|x_{N-1}, \dots, x_2, x_1) = \\ &= \prod_{i=1}^N p(x_i|x_{i-1}, \dots, x_2, x_1) \end{aligned}$$

$$p(x_1, x_2) = p(x_1)p(x_2|x_1) = p(x_2)p(x_1|x_2) \xrightarrow[\substack{\text{independent} \\ X_1 \perp X_2}]{} p(x_1)p(x_2)$$

$$p_{\mathbf{x}}(x_1, x_2, \dots, x_N) = p_{x_1}(x_1)p_{x_2}(x_2)p_{x_3}(x_3)\dots p_{x_N}(x_N) = \prod_{i=1}^N p_{x_i}(x_i)$$

For the **independent and identically distributed (i.i.d.)** random variables

$$p_{\mathbf{x}}(x_1, x_2, \dots, x_N) = p_x(x_1)p_x(x_2)p_x(x_3)\dots p_x(x_N) = \prod_{i=1}^N p_x(x_i)$$

# Random vectors: independence

---

- Practical application of independence

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{x}}(x_1, x_2, \dots, x_N)$$

$\mathbf{X} =$

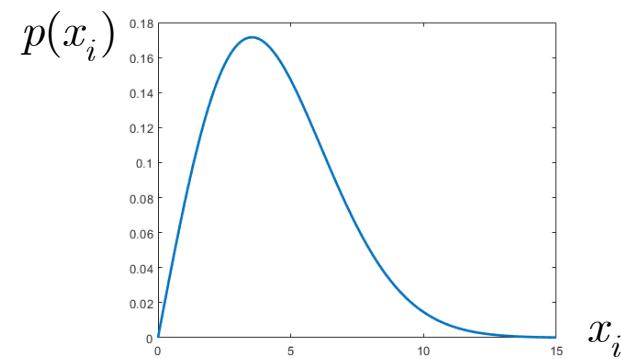


$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

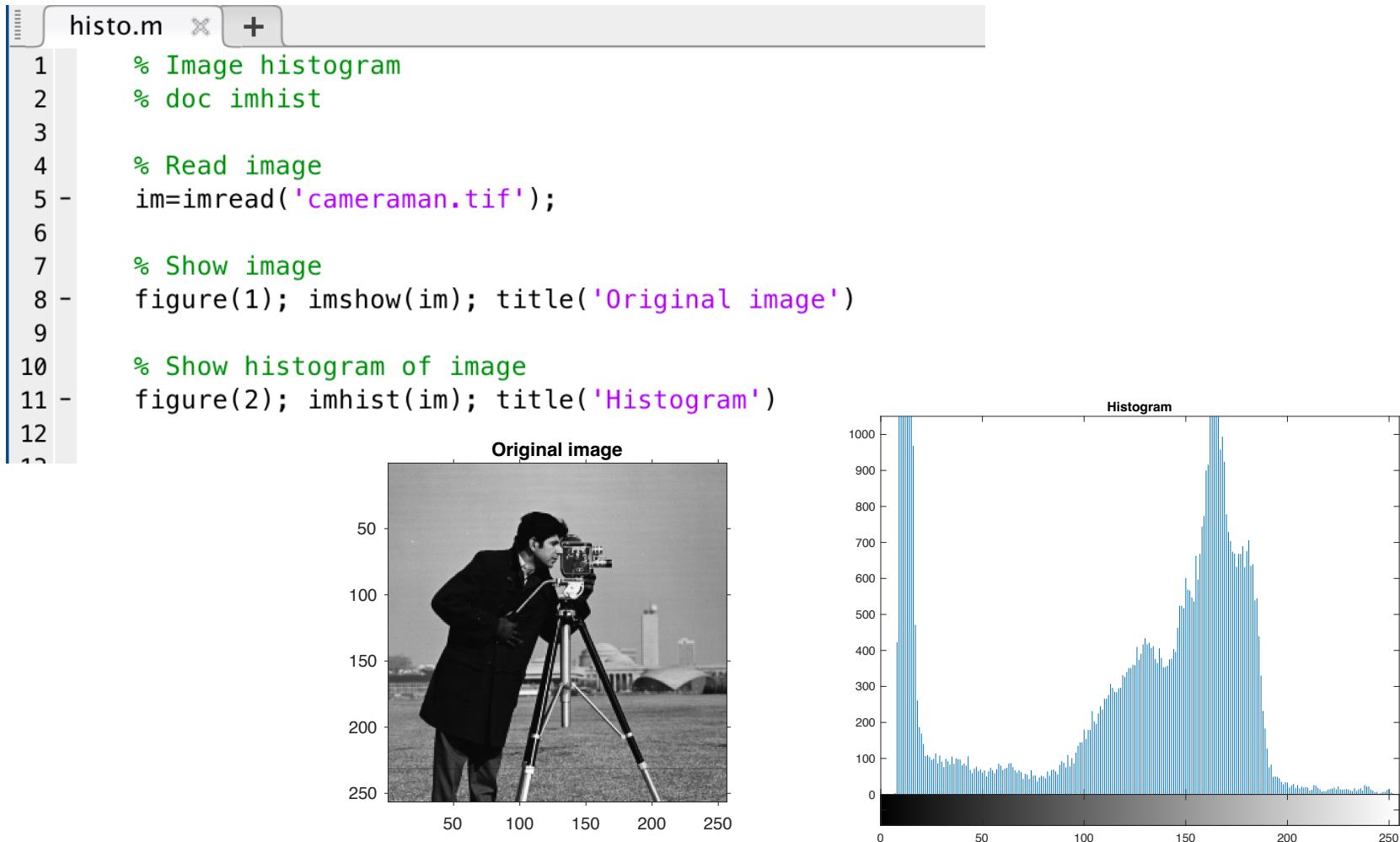
Simplification

$$p(x_1, x_2, \dots, x_N) = p(x_1)p(x_2)p(x_3)\dots p(x_N)$$

$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix} \sim p(x_1) \sim p(x_2) \sim p(x_N)$$



# Matlab example of histogram of image pixels

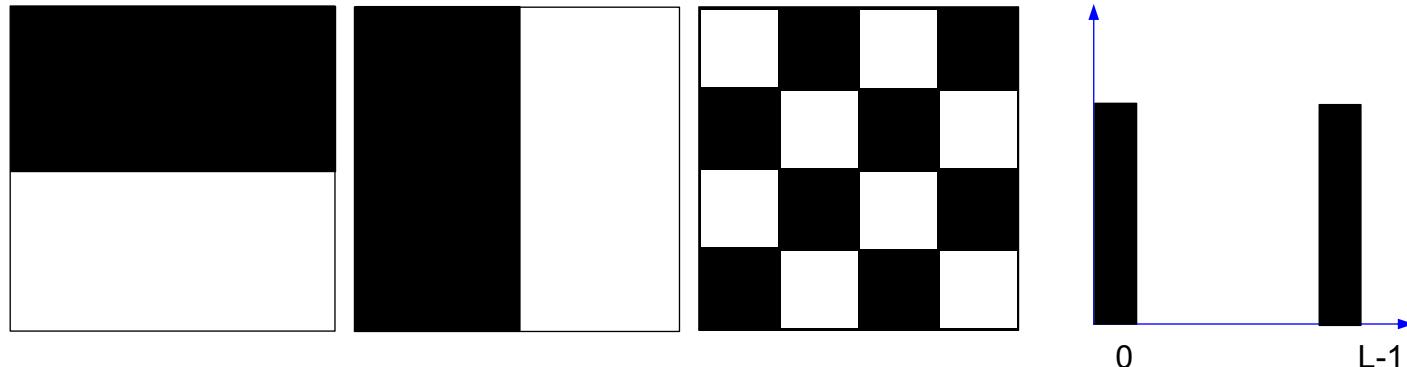


# Remarks about histogram

---

Is histogram informative? Yes

Is histogram unique? Many images might have the same histogram!



**Important: histogram is not a pdf or pmf!**

# Random vectors: moments

---

- If random variables  $X$  and  $Y$  are independent ( $X \perp Y$ )

$$E[XY] = E[X]E[Y]$$

$$E[X|Y=y] = E[X]$$

$$E[Y|X=x] = E[Y]$$

$$\text{Cov}[X,Y] = 0, \rho_{X,Y} = 0$$

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$$

- Proof:

- For two independent random variables  $p_{X,Y}(x,y) = p_X(x)p_Y(y)$

$$E[XY] = E[X]E[Y]$$

$$E[XY] = \sum_{x \in S_X} \sum_{y \in S_Y} xy p_X(x)p_Y(y) = \sum_{x \in S_X} xp_X(x) \sum_{y \in S_Y} yp_Y(y) = E[X]E[Y]$$

# Indépendance entre v.a.

---

$$Cov[X, Y] = 0, \rho_{X,Y} = 0$$

$$Cov[X, Y] = E[XY] - \mu_X\mu_Y = E[X]E[Y] - \mu_X\mu_Y = 0$$

$$\rho_{XY} = \frac{Cov[X, Y]}{\sqrt{Var[X] Var[Y]}} = 0$$

$$Var[X + Y] = Var[X] + Var[Y]$$

$$Var[X \pm Y] = Var[X] + Var[Y] \pm 2\underbrace{Cov[X, Y]}_{=0}$$

# Random vectors: Markovianity

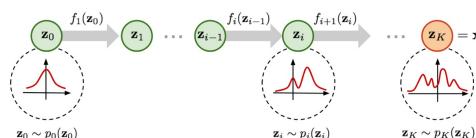
- The **first order** Markov model
  - assumes the dependence on **one element** from past

$$p(x_i | x_{i-1}, x_{i-2}, \dots, x_1) \underset{1st \text{ order Markov model}}{\Rightarrow} p(x_i | x_{i-1})$$

$$p(x_1, x_2, \dots, x_N) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1)\dots p(x_N|x_{N-1}, \dots, x_2, x_1) = \\ \underset{1st \text{ order Markov model}}{\Rightarrow} p(x_1)p(x_2|x_1)p(x_3|x_2)\dots p(x_N|x_{N-1})$$

$$p(x_1, x_2, \dots, x_N) \underset{1st \text{ order Markov model}}{\Rightarrow} p(x_1) \prod_{i=2}^N p(x_i | x_{i-1})$$

The modeling of a complex pmf/pdf is reduced to the definition of pair-wise conditional distributions



HAMILTONIAN GENERATIVE NETWORKS

# Random vectors: Markovianity

- The **second order** Markov model
  - assumes the dependence on **two elements** from past

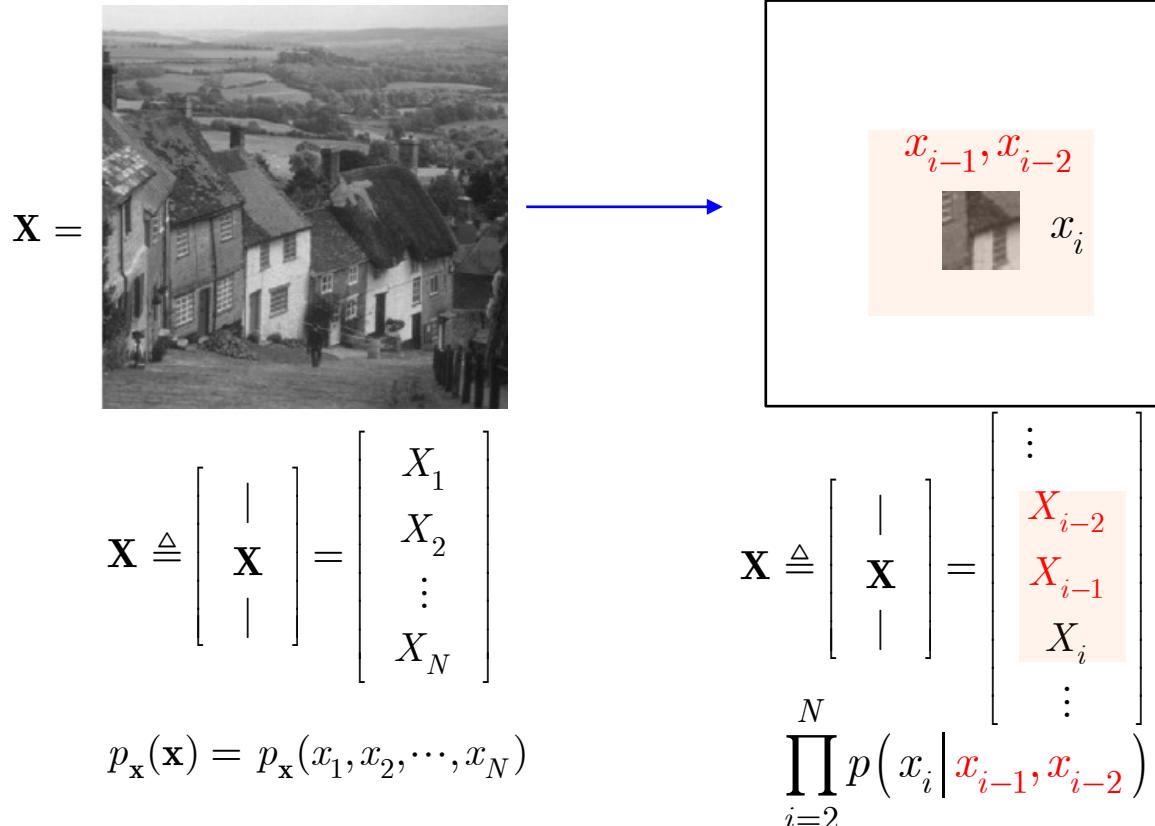
$$p(x_i | x_{i-1}, x_{i-2}, \dots, x_1) \xrightarrow{\text{2nd order Markov model}} p(x_i | x_{i-1}, x_{i-2})$$

$$\begin{aligned} p(x_1, x_2, \dots, x_N) &= p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \dots p(x_N | x_{N-1}, \dots, x_2, x_1) = \\ &\xrightarrow{\text{2nd order Markov model}} p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) \dots p(x_N | x_{N-1}, x_{N-2}) \end{aligned}$$

$$p(x_1, x_2, \dots, x_N) \xrightarrow{\text{2nd order Markov model}} p(x_1) \prod_{i=2}^N p(x_i | x_{i-1}, x_{i-2})$$

# Random vectors: Markovianity

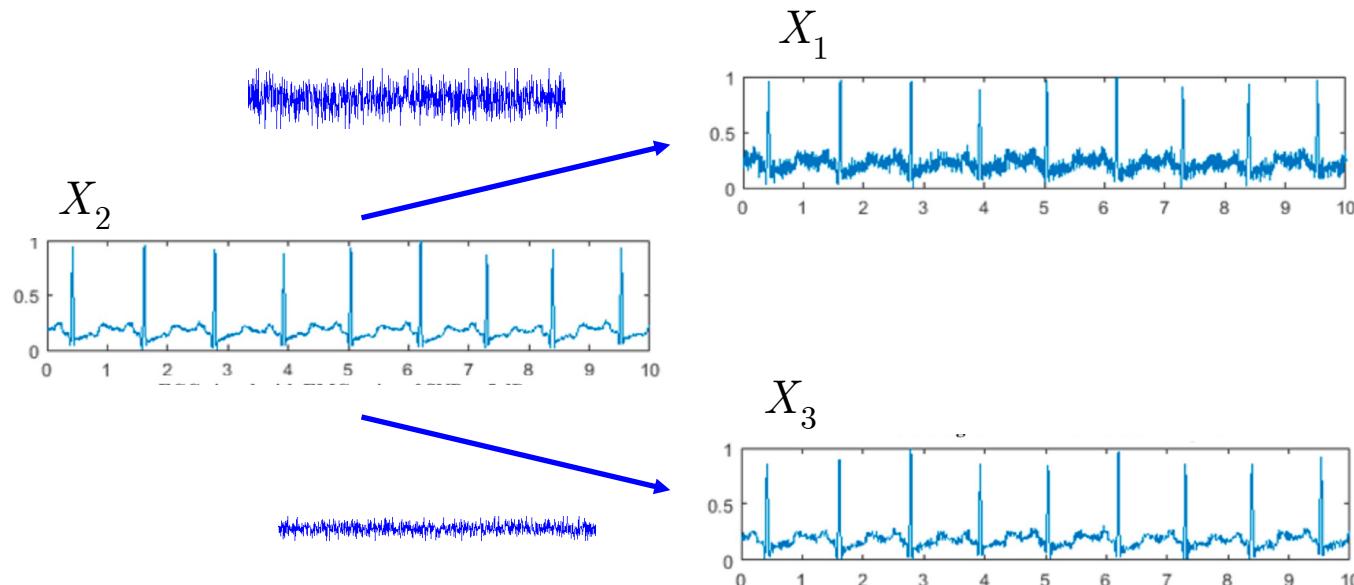
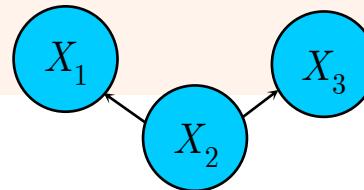
- Practical application application of Markovianity for image processing



# Random vectors: conditional independence

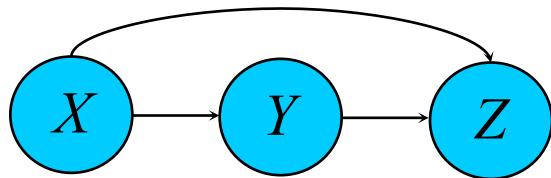
- Random variables  $X_1$  and  $X_3$  are independent given  $X_2$

$$p(x_1, x_2, x_3) = p(x_2) \underbrace{p(x_1, x_3 | x_2)}_{= p(x_1 | x_2) p(x_3 | \textcolor{red}{x}_1, x_2)} = p(x_2) p(x_1 | x_2) p(x_3 | x_2)$$

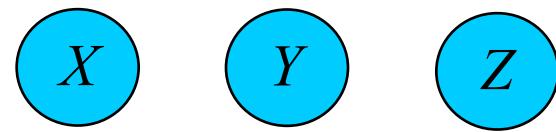


# Summary

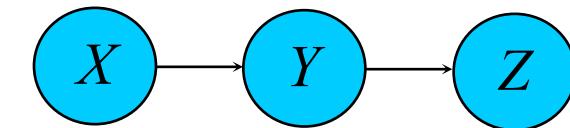
---



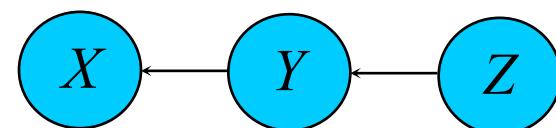
$$p(x, y, z) = p(x)p(y|x)p(z|x, y)$$



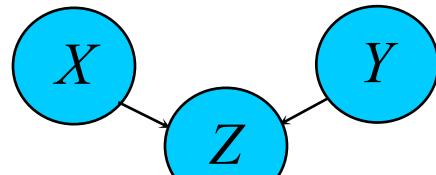
$$p(x, y, z) = p(x)p(y)p(z)$$



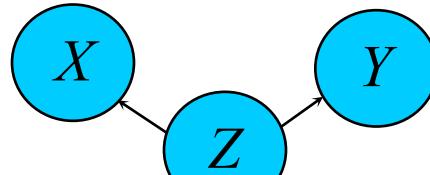
$$p(x, y, z) = p(x)p(y|x)p(z|y)$$



$$p(x, y, z) = p(z)p(y|z)p(x|y)$$



$$p(x, y, z) = p(z|x, y)p(x)p(y)$$



$$p(x, y, z) = p(z)p(x|z)p(y|z)$$

# Random vectors: marginalization

---

- Marginalization:
  - Given  $p(x_1, x_2)$ , one can find marginal distributions by

$$p(x_1) = \sum_{x_2 \in S_{X_2}} p(x_1, x_2) = \sum_{x_2 \in S_{X_2}} p(x_1) p(x_2 | x_1)$$

$$p(x_2) = \sum_{x_1 \in S_{X_1}} p(x_1, x_2) = \sum_{x_1 \in S_{X_1}} p(x_2) p(x_1 | x_2)$$

# Transformation of random variables

---

- General linear transform:
  - Random scalar variables
  - Random vectors
- Sum of random variables
  - Sum of independent random variables
  - Concentration inequalities
  - Weak Law of Large Numbers and CLT
  - Application to vector norms (see later after norms)
- Transformation of Gaussian random vectors
  - Mean and covariance matrix after transformation
  - Pdf after transformation
  - Properties of transform
    - Decorrelation
    - Whitening

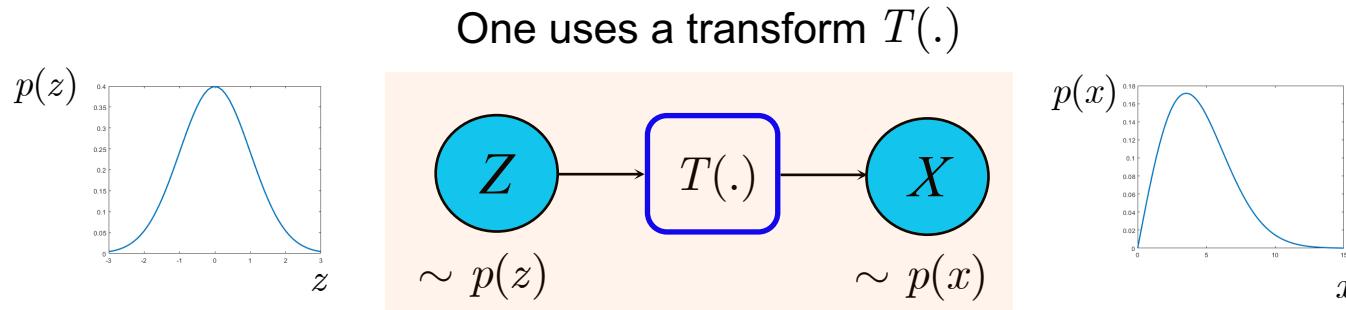
# Transformation of random variables

---

- Why are we interested in transformation of random variables/vectors?
  
- The transformation is a core operation for:
  - Data analysis and visualization (looking for interesting features in data)
  - Fast data indexing and searching (lower dim data)
  - Data processing
    - Clustering (grouping of similar data), denoising (removal of noise), restoration (recovery from distortion), classification (assigning to the predefined class labels), etc.
  - Data compression (more efficient/compact data representation)
  - Data transmission (robust data encoding)
  - Data generation (generation of complex distributions from simple ones)

# Transformation of random variables

- **Goal:** to transform a random variable  $Z \sim p(z)$  into another random variable  $X \sim p(x)$ 
  - It is assumed that  $Z$  follows some simple distribution that can be easily generated and analyzed (uniform or Gaussian)
  - While the distribution of  $X$  is complex (asymmetric, multimodal, etc.).
  - SOTA generative models (VAE, GAN, FLOW)
- **Scalar-to-scalar:**  $T : \mathbb{R} \rightarrow \mathbb{R}$



Note: formally we need to use  $f_Z(z)$  instead of  $p(z)$  for continuous variables

# Transformation of random variables

---

- Properties of transform  $T(\cdot)$ :
  - Invertible  $T^{-1}$
  - Both  $T$  and  $T^{-1}$  are differentiable
- Generate  $Z \sim p(z)$
- Consider one-to-one mapping transform:  $x = T(z)$ 
  - Recalling  $\int p(z)dz = \int p(x)dx = 1$
  - Preservation of “volume”:
    - Probability is invariant under the change of variables  $|p(x)dx| = |p(z)dz|$

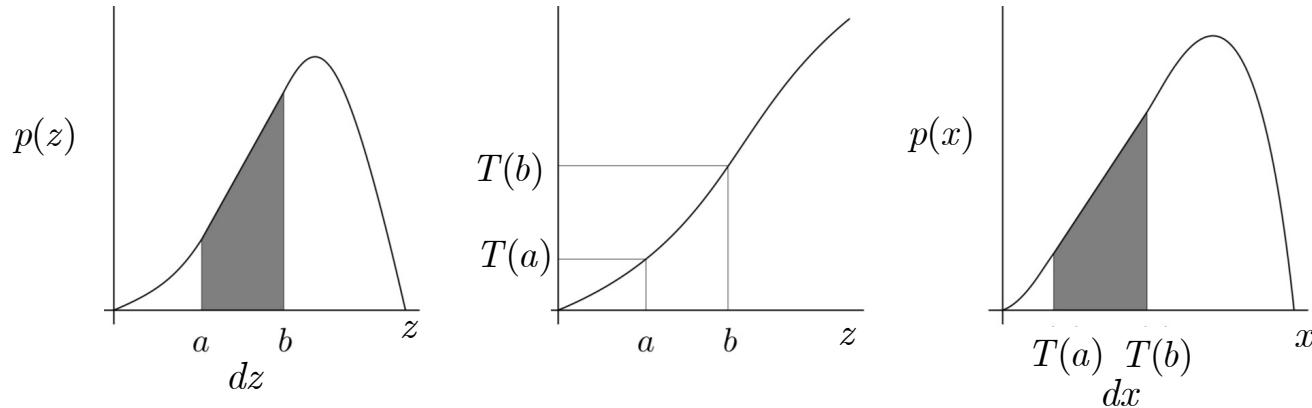
$$p(x) = p(z) \left| \frac{dz}{dx} \right| = p(T^{-1}(x)) \left| \frac{dT^{-1}(x)}{dx} \right|$$

$$\begin{aligned} x &= T(z) \\ z &= T^{-1}(x) \end{aligned}$$

# Transformation of random variables

---

- Example (more details)



$$\Pr[a \leq Z < b] = \int_a^b p(z) dz = \Pr[T(a) \leq X < T(b)] = \int_{T(a)}^{T(b)} p(x) dx$$

$$|p(z)dz| = |p(x)dx|$$

$$p(x) = p(z) \left| \frac{dz}{dx} \right|$$

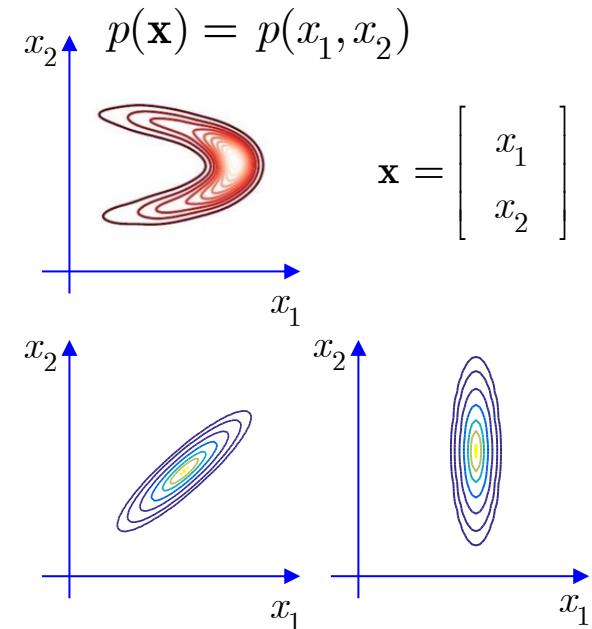
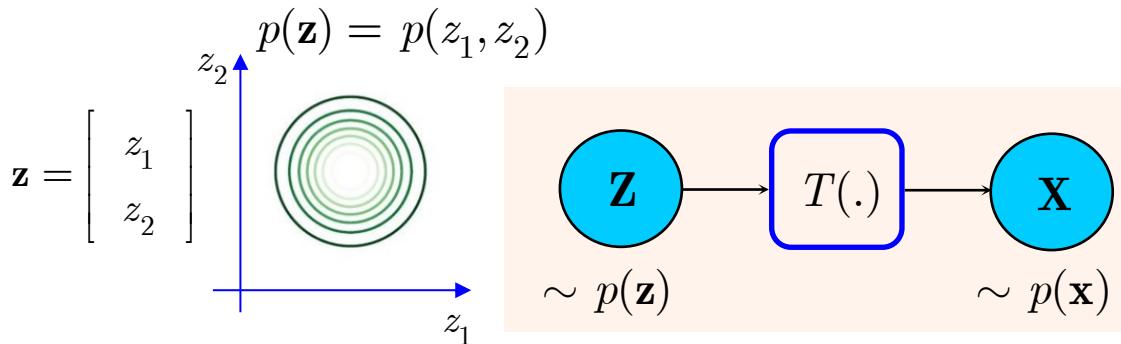
# Transformation of random variables

- **Vector-to-vector:**  $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$

$$\mathbf{x} = T(\mathbf{z})$$

$$\mathbf{z} = T^{-1}(\mathbf{x})$$

$$\begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = T \begin{pmatrix} | \\ \mathbf{z} \\ | \end{pmatrix} = \begin{bmatrix} T_1(\mathbf{z}) \\ T_2(\mathbf{z}) \\ \vdots \\ T_N(\mathbf{z}) \end{bmatrix}$$



# Transformation of random vectors

---

- **Vector-to-vector:**  $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$

$$\mathbf{x} = T(\mathbf{z})$$

$\mathbf{z} = T^{-1}(\mathbf{x})$ , i.e., the transform is invertible

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) \left| \det J_T(\mathbf{z}) \right|^{-1}$$

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(T^{-1}(\mathbf{x}))}{\left| \det J_T(\mathbf{z}) \right|}$$

The Jacobian matrix:

$$J_T(\mathbf{z}) = \begin{bmatrix} \frac{\partial T_1(\mathbf{z})}{\partial z_1} & \dots & \frac{\partial T_1(\mathbf{z})}{\partial z_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_N(\mathbf{z})}{\partial z_1} & \dots & \frac{\partial T_N(\mathbf{z})}{\partial z_N} \end{bmatrix} \quad J_T(\mathbf{x}) = \begin{bmatrix} \frac{\partial T_1^{-1}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial T_1^{-1}(\mathbf{x})}{\partial x_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_N^{-1}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial T_N^{-1}(\mathbf{x})}{\partial x_N} \end{bmatrix}$$

# Transformation of random variables

---

- Example:  $T : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\mathbf{x} = \begin{bmatrix} T_1(\mathbf{z}) \\ T_2(\mathbf{z}) \end{bmatrix} = \begin{bmatrix} T_1(z_1, z_2) \\ T_2(z_1, z_2) \end{bmatrix} = \begin{bmatrix} z_1^2 z_2 \\ 5z_1 + \sin(z_2) \end{bmatrix}$$

The Jacobian matrix:

$$J_T(\mathbf{z}) = \begin{bmatrix} \frac{\partial T_1(\mathbf{z})}{\partial z_1} & \frac{\partial T_1(\mathbf{z})}{\partial z_2} \\ \frac{\partial T_2(\mathbf{z})}{\partial z_1} & \frac{\partial T_2(\mathbf{z})}{\partial z_2} \end{bmatrix} = \begin{bmatrix} 2z_1 z_2 & z_1^2 \\ 5 & \cos(z_2) \end{bmatrix}$$

The Jacobian determinant:

$$|\det J_T(\mathbf{z})| = 2z_1 z_2 \cos(z_2) - 5z_1^2$$

# Linear transformation of random vectors

---

- **Vector-to-vector:**  $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$

$$\mathbf{x} = T(\mathbf{z}) = \mathbf{A}\mathbf{z}$$

$\mathbf{z} = T^{-1}(\mathbf{x}) = \mathbf{A}^{-1}\mathbf{x}$ , i.e., the transform is invertible

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(T^{-1}(\mathbf{x}))}{|\det J_T(\mathbf{z})|} = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det J_T(\mathbf{z})|}$$

The Jacobian matrix:

$$J_T(\mathbf{z}) = \begin{bmatrix} \frac{\partial T_1(\mathbf{z})}{\partial z_1} & \dots & \frac{\partial T_1(\mathbf{z})}{\partial z_N} \\ \vdots & \ddots & \vdots \\ \frac{\partial T_N(\mathbf{z})}{\partial z_1} & \dots & \frac{\partial T_N(\mathbf{z})}{\partial z_N} \end{bmatrix}$$

# Linear transformation of random vectors

- **Vector-to-vector:**  $T : \mathbb{R}^N \rightarrow \mathbb{R}^N$

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(T^{-1}(\mathbf{x}))}{|\det J_T(\mathbf{z})|} = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det J_T(\mathbf{z})|}$$

The Jacobian matrix:

$$\mathbf{x} = \mathbf{A}\mathbf{z} = \begin{bmatrix} a_{11} & \cdots & a_{1N} \\ \vdots & \ddots & \vdots \\ a_{N1} & \cdots & a_{NN} \end{bmatrix} \begin{bmatrix} z_1 \\ \vdots \\ z_N \end{bmatrix}$$

$$x_j = T_j(\mathbf{z}) = a_{j1}z_1 + \cdots + a_{ji}z_i + \cdots + a_{jN}z_N \quad \frac{\partial T_j(\mathbf{z})}{\partial z_i} = \frac{\partial x_j}{\partial z_i} = a_{ji} \quad \Rightarrow J_T(\mathbf{z}) = \mathbf{A}$$

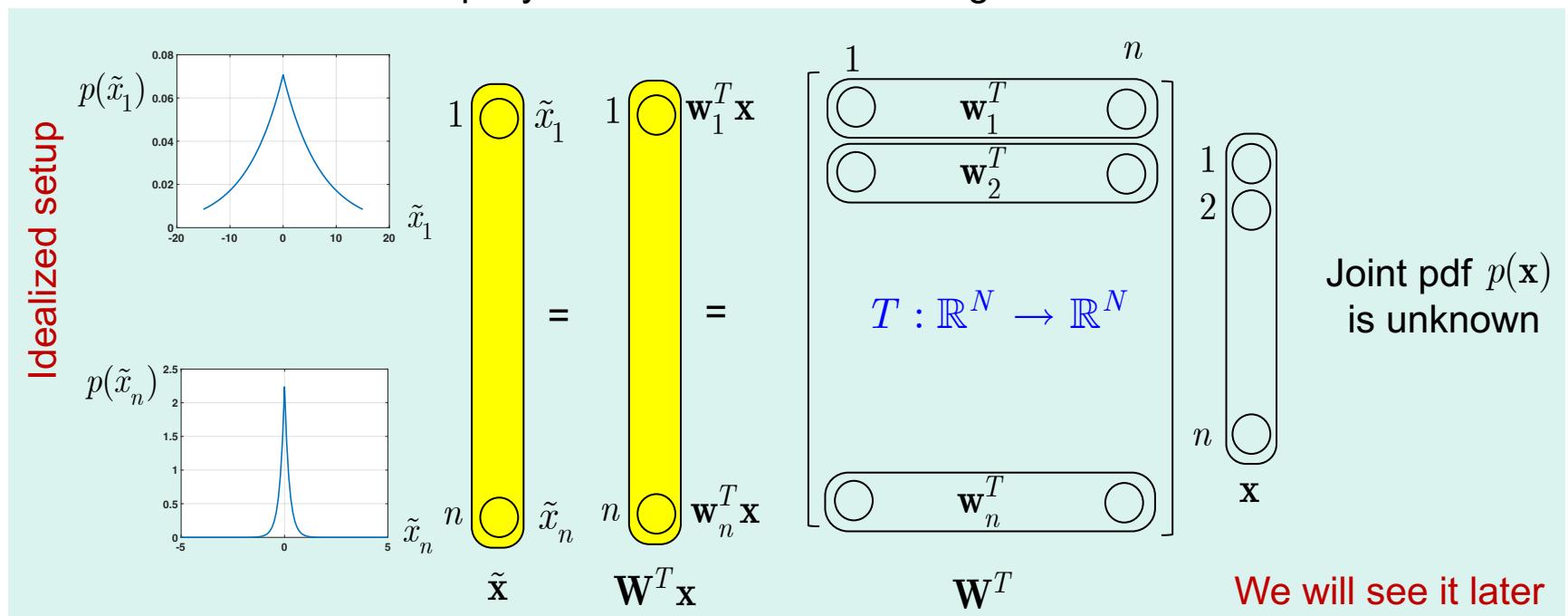
$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(T^{-1}(\mathbf{x}))}{|\det J_T(\mathbf{z})|} = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det \mathbf{A}|}$$

# Transformation for decorrelation and compaction

- **Main idea:** assume  $\mathbf{x}$  can be transformed to some domain where all pixels will be independent (ideally; just a wish), i.e.:

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_n) = p(x_1)p(x_2|x_1)\cdots p(x_n|x_{n-1}, \dots, x_1) \Rightarrow p(\tilde{\mathbf{x}}) = p(\tilde{x}_1)p(\tilde{x}_2)\cdots p(\tilde{x}_n)$$

- such that we can uniquely transform back to the original vector



# Transformation for decorrelation and compaction

---

- The key element of transformation is a projection

$$\mathbf{w}_i^T \mathbf{x} = \sum_{j=1}^N w_{i,j} x_j = w_{i,1} x_1 + w_{i,2} x_2 + \cdots + w_{i,N} x_N$$

- Assuming  $\left\{ w_{i,j} \right\}_{j=1}^N = 1$ , we have just a sum

$$s_N(x) = x_1 + x_2 + \cdots + x_N$$

- Consider all variables to be random ones, we obtain a random value of sum

$$S_N(X) = X_1 + X_2 + \cdots + X_N$$

- Since the sum of random variables plays a fundamental role for many concentration properties, we will consider it next.

# Transformation of random variables

---

- General linear transform:
  - Random scalar variables
  - Random vectors
- Sum of random variables
  - Sum of independent random variables
  - Concentration inequalities
  - Weak Law of Large Numbers and CLT
  - Application to vector norms (see later after norms)
- Transformation of Gaussian random vectors
  - Mean and covariance matrix after transformation
  - Pdf after transformation
  - Properties of transform
    - Decorrelation
    - Whitening

# Sum of random variables: pdf and moments

---

- Sum of random variables  $S_N(X) = X_1 + X_2 + \dots + X_N$

- It is also often of interest to compute a sampling mean

$$M_N(X) = \frac{S_N(X)}{N} = \frac{X_1 + X_2 + \dots + X_N}{N}$$

- The statistical behavior of  $S_N(X)$  is characterized by:

- **Pdf**  $f_{X_1+\dots+X_N}(x)$
  - If it is difficult to compute pdf, one computes the first two moments
    - **Expected value of sum**

$$E[S_N(X)] = E[X_1 + \dots + X_N] = E[X_1] + \dots + E[X_N]$$

- **Variance of sum**

$$Var[S_N(X)] = \sum_{i=1}^N Var[X_i] + \sum_{i=1}^N \sum_{j \neq i} Cov[X_i, X_j]$$

# Sum of random variables: pdf and moments

---

- Sum of **uncorrelated** random variables

- Variance

$$Cov[X_i, X_j] = 0$$

$$Var[S_N(X)] = \sum_{i=1}^N Var[X_i] + \sum_{i=1}^N \sum_{j=1}^N \underbrace{Cov[X_i, X_j]}_{=0} = \sum_{i=1}^N Var[X_i]$$

- Sum of **independent** random variables

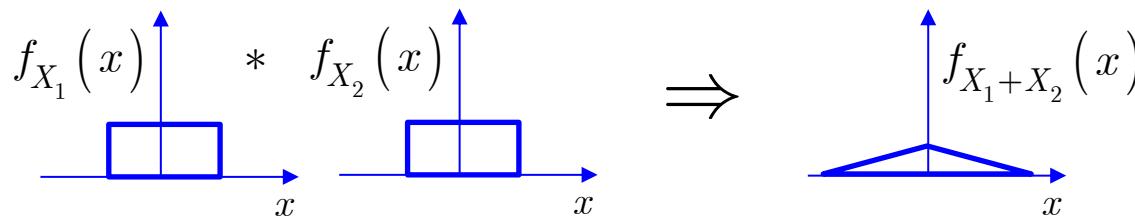
$$f_{X_1+\dots+X_N}(x) = (f_{X_1}(x) * \dots * f_{X_N}(x))(x)$$

where  $f_{X_1+X_2}(x) = \int_{-\infty}^{+\infty} f_{X_1}(u) f_{X_2}(x-u) du = (f_{X_1}(x) * f_{X_2}(x))(x)$

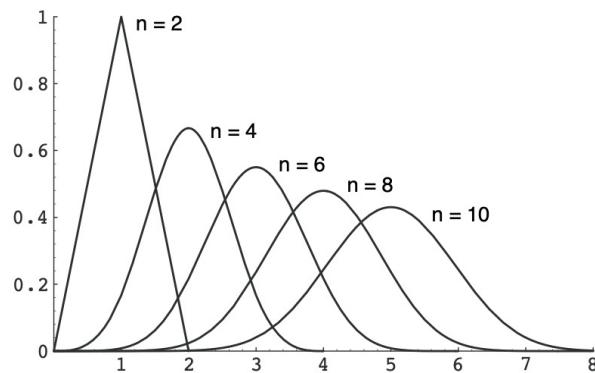
**convolution**

# Sum of random variables: convergence to Gaussian

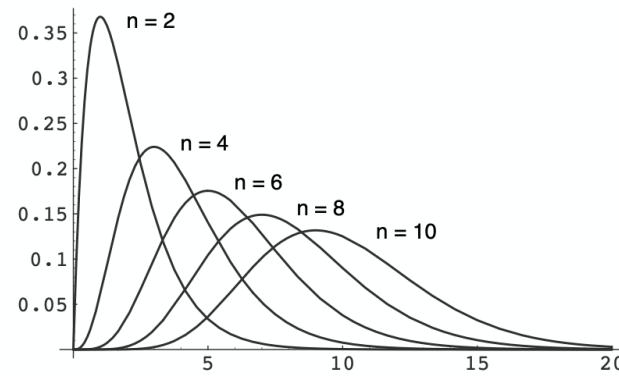
- Sum of **independent** random variables
  - Example: sum of two uniform independent random variables



Convolution of  $N = n$  uniform pdfs



Convolution of  $N = n$  exponential pdfs



- The pdf of sum converges to Gaussian pdf

# Sample mean

---

- Consider a sample mean of  $N$  independent and identically distributed random variables  $X_1, X_2, \dots, X_N$  with  $E[X_i] = \mu_X$  and  $\text{Var}[X_i] = \sigma_X^2$
- The **sample mean** is

$$M_N(X) = \frac{S_N(X)}{N} = \frac{X_1 + X_2 + \dots + X_N}{N} = \frac{1}{N} \sum_{i=1}^N X_i$$

- The moments of **sample mean**
  - **Expected value:**  $E[M_N(X)] = E[X] = \mu_X$
  - **Variance:**  $\text{Var}[M_N(X)] = \frac{\text{Var}[X]}{N} = \frac{\sigma_X^2}{N}$

# Sample mean: moments

---

- **Proof**
- **Expected value:** sampling mean is an **unbiased estimator**

$$E[M_N(X)] = \frac{1}{N}(E[X_1] + \cdots + E[X_N]) = \frac{1}{N}(E[X] + \cdots + E[X]) = E[X] = \mu_X$$

- **Variance:** suppose that  $Var[X_i] = \sigma_X^2$

$$\begin{aligned} Var[M_N(X)] &= \frac{1}{N^2}(Var[X_1] + \cdots + Var[X_N]) \\ &= \frac{1}{N^2}(Var[X] + \cdots + Var[X]) = \frac{1}{N}Var[X] = \frac{\sigma_X^2}{N} \end{aligned}$$

The variance of sample mean approaches zero as the number of samples increases  
(a.k.a. **efficient estimator**)

# Weak Law of Large Numbers (WLLN)

- Let  $X_1, X_2, \dots, X_N$  be a sequence of i.i.d. random variables with finite mean  $E[X_i] = \mu_X$  and variance  $\text{Var}[X_i] = \sigma_X^2$
- For any  $\varepsilon > 0$ , **Chebyshev inequality** states:

$$(a) \Pr[|M_N(X) - \mu_X| \geq \varepsilon] \leq \frac{\sigma_X^2}{N\varepsilon^2}$$

$$(b) \Pr[|M_N(X) - \mu_X| < \varepsilon] \geq 1 - \frac{\sigma_X^2}{N\varepsilon^2}$$

- Weak Law of Large Numbers**

$$(a) \lim_{N \rightarrow \infty} \Pr[|M_N(X) - \mu_X| \geq \varepsilon] = 0$$

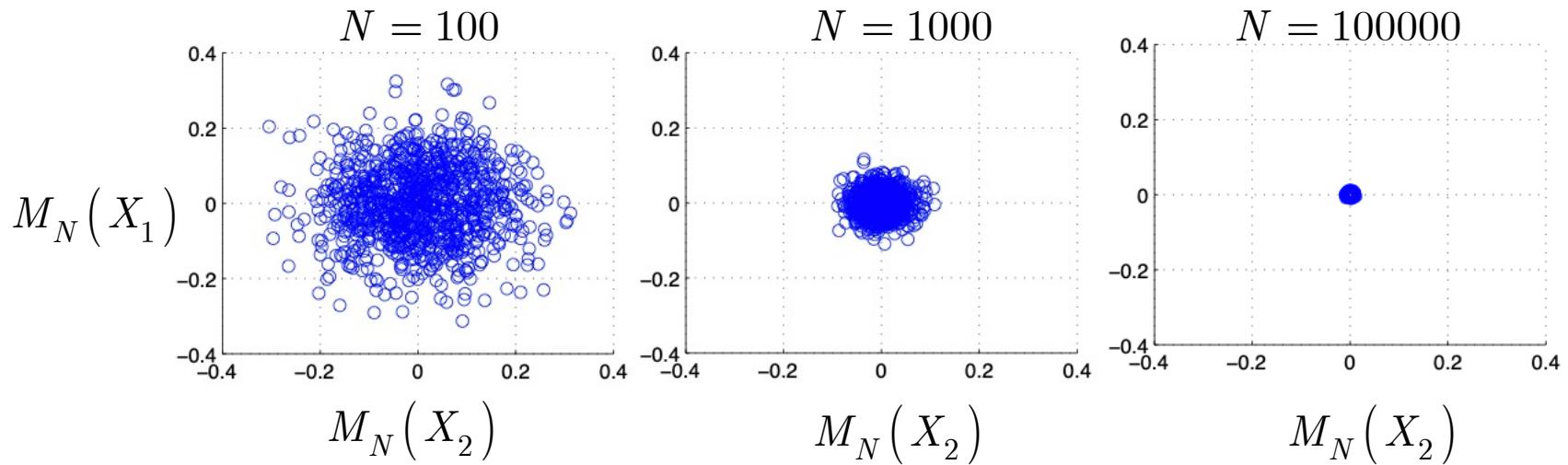
$$(b) \lim_{N \rightarrow \infty} \Pr[|M_N(X) - \mu_X| < \varepsilon] = 1$$

# Weak Law of Large Numbers (WLLN)

---

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad X_1 \perp X_2 \sim f_X(x) = \mathcal{N}(0, \sigma_X^2)$$

$$M_N(X_1) = \frac{S_N(X_1)}{N} = \frac{1}{N} \sum_{i=1}^N X_{1,i} \xrightarrow[N \rightarrow \infty]{} 0 \quad M_N(X_2) = \frac{S_N(X_2)}{N} = \frac{1}{N} \sum_{i=1}^N X_{2,i} \xrightarrow[N \rightarrow \infty]{} 0$$

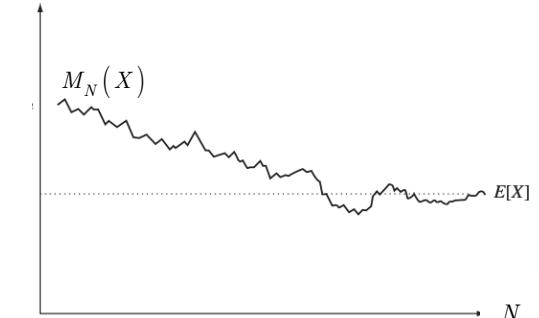


# Strong Law of Large Numbers (SLLN)

---

- Let  $X_1, X_2, \dots, X_N$  be a sequence of i.i.d. random variables with finite mean  $E[X_i] = \mu_X$  and variance  $\text{Var}[X_i] = \sigma_X^2$  (some bounded)
- In this case

$$\Pr\left[\lim_{N \rightarrow \infty} M_N(X) = \mu_X\right] = 1$$



- $M_K(X)$  is the sample mean computed using  $K$  variables
- With probability 1, every sample mean calculation will eventually approach and stay close to  $E[X_i] = \mu_X$
- The strong law implies the weak law

# Central Limit Theorem (CLT)

- Let  $X_1, X_2, \dots, X_N$  be a sequence of i.i.d. random variables with finite mean  $E[X_i] = \mu_X$  and variance  $\text{Var}[X_i] = \sigma_X^2$
- Let  $S_N(X)$  be the sum:

$$S_N(X) = X_1 + X_2 + \dots + X_N$$

- Define the centered and normalized variable

$$Z_N(X) = \frac{S_N(X) - N\mu_X}{\sigma_X \sqrt{N}}$$

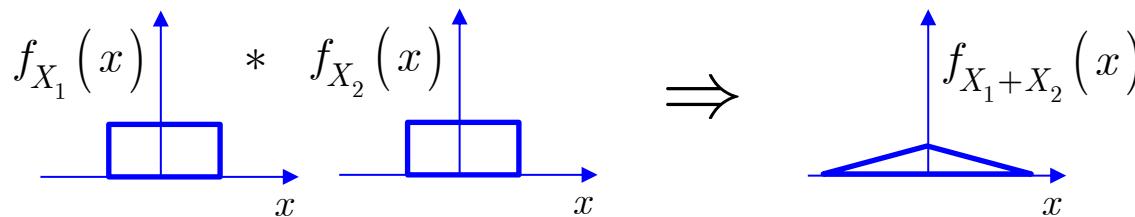
- then

$$\lim_{N \rightarrow \infty} \Pr[Z_N(X) \leq z] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx$$

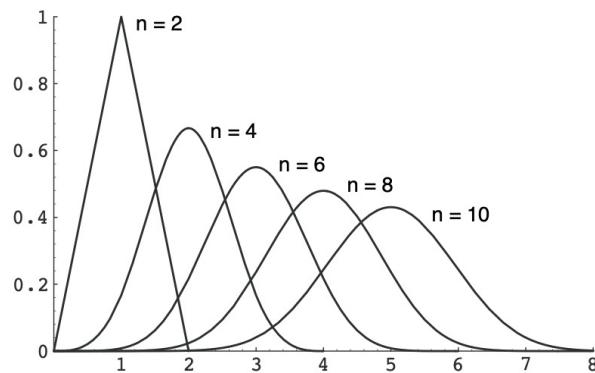
as  $N$  becomes large, the CDF of normalized and centered  $Z_N(X)$  approaches Gaussian pdf

# Sum of random variables: convergence to Gaussian

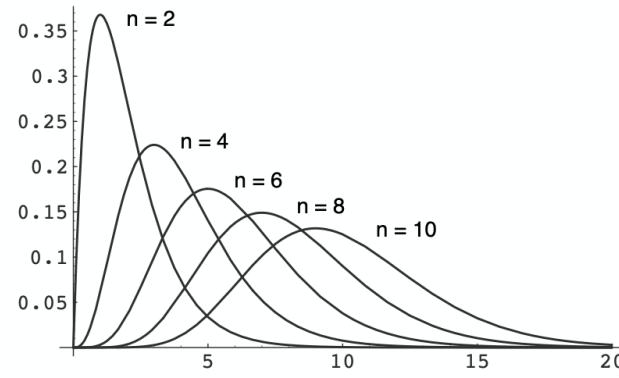
- Sum of **independent** random variables
  - Example: sum of two uniform independent random variables



Convolution of  $N = n$  uniform pdfs



Convolution of  $N = n$  exponential pdfs



- The pdf of sum converges to Gaussian pdf

# Similarity between vectors

---

- We have seen the “similarity” between random variables expressed by covariance (via the second order moments) or independence (the level of pmf/pdf)
- Our next objectives:
  - To introduce a measure of similarity for deterministic vectors via a notion of metric and define:
    - Distance, inner product and cosine similarity
  - To extend it to random vectors:
    - Expected value of distance and inner product
    - Concentration properties (weak law of large numbers)
  - To extend the notion of similarity to distributions:
    - Divergence: KLD and f-div (Theme 2)
    - Mutual Information (Theme 2)

# Vector norms

- Vector norms  $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$

$$\mathbf{x} \in \mathbb{R}^N$$

$$\|\mathbf{x}\|_p$$

$\ell_p$  - norm or  $p$ - norm

$$\|\mathbf{x}\|_p = \left( |x_1|^p + |x_2|^p + \dots + |x_N|^p \right)^{\frac{1}{p}} = \left( \sum_{i=1}^N |x_i|^p \right)^{\frac{1}{p}}$$

$$\|\mathbf{x}\|_2$$

$\ell_2$  - norm (Euclidian)

$$\|\mathbf{x}\|_2 = \left( |x_1|^2 + |x_2|^2 + \dots + |x_N|^2 \right)^{\frac{1}{2}} = \left( \sum_{i=1}^N |x_i|^2 \right)^{\frac{1}{2}}$$

$$\|\mathbf{x}\|_1$$

$\ell_1$  - norm

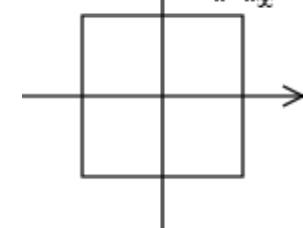
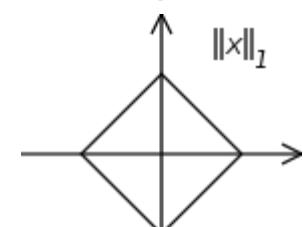
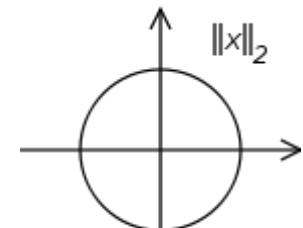
$$\|\mathbf{x}\|_1 = (|x_1| + |x_2| + \dots + |x_N|) = \left( \sum_{i=1}^N |x_i| \right)$$

$$\|\mathbf{x}\|_\infty$$

$\ell_\infty$  - norm or maximum norm

$$\|\mathbf{x}\|_\infty = \max \{|x_1|, |x_2|, \dots, |x_N|\}$$

Unit-norm



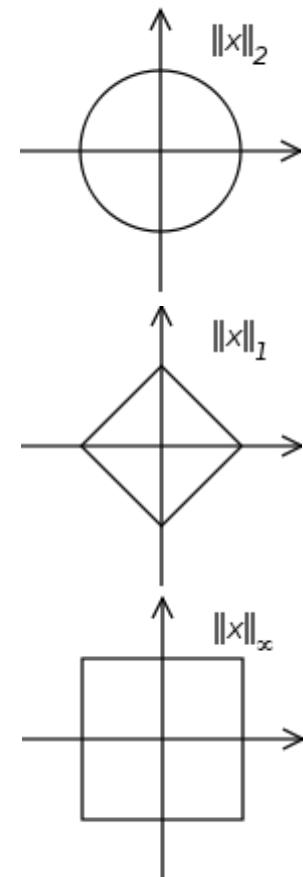
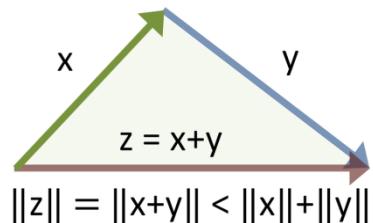
# Vector norms

- Vector norms  $\mathbf{x} = (x_1, x_2, \dots, x_N)^T$

$$\mathbf{x} \in \mathbb{R}^N$$

For all  $p \geq 1$ , the **p-norms and maximum norm** as defined above satisfy the properties of a "length function" (or norm), which are that:

- only the zero vector has zero length,
- the length of the vector is positive homogeneous with respect to multiplication by a scalar ( $f(\alpha\mathbf{v}) = \alpha^k f(\mathbf{v})$ )
- the length of the sum of two vectors is not larger than the sum of lengths of the vectors (triangle inequality).

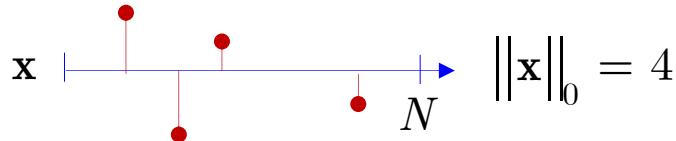


# Vector "norm" for sparse vectors

---

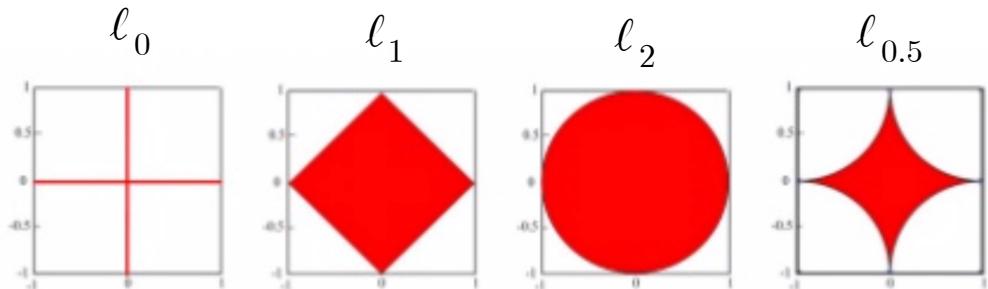
- The **number of nonzeros** is called the  $\ell_0$ -"norm" of  $\mathbf{x}$

$$\|\mathbf{x}\|_0 \doteq \#\{i : x_i \neq 0\}$$



$$\|\mathbf{x}\|_p = \left( \sum_{i=1}^N |x_i|^p \right)^{1/p}$$

$$\|\mathbf{x}\|_0 = \lim_{p \rightarrow 0} \|\mathbf{x}\|_p^p$$



# Distances

---

## Distances and similarity measures

A **metric**  $d$  is a function (mapping)

$$d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$$

that satisfies:

1.  $d(x, y) \geq 0$  (*non-negative*)
2.  $d(x, y) = 0$ , iff  $x = y$  (*identity of indiscernibles*)
3.  $d(x, y) = d(y, x)$  (*symmetry*)
4.  $d(x, z) \leq d(x, y) + d(y, z)$  (*triangle inequality*)

# Distances

---

## Distances and similarity measures

$$\|\mathbf{x} - \mathbf{y}\|_2$$

Euclidian distance between two vectors

$$d_2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2 = \left( \sum_{i=1}^N (x_i - y_i)^2 \right)^{\frac{1}{2}}$$

$$\|\mathbf{x} - \mathbf{y}\|_1$$

$\ell_1$  - distance between two vectors

$$d_1(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^N |x_i - y_i|$$

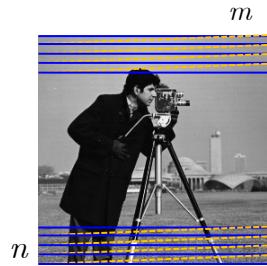
$$\mathbf{x} \oplus \mathbf{y}$$

Hamming distance between two vectors

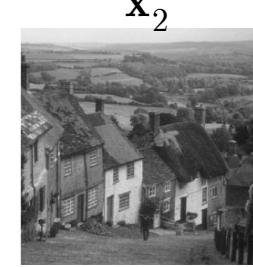
$$d^H(\mathbf{x}, \mathbf{y}) = \mathbf{x} \oplus \mathbf{y} = \sum_{i=1}^N 1\{x_i, y_i\}$$

with indicator function  $1\{x_i, y_i\} = \begin{cases} 0, & \text{if } x_i = y_i \\ 1, & \text{if } x_i \neq y_i \end{cases}$

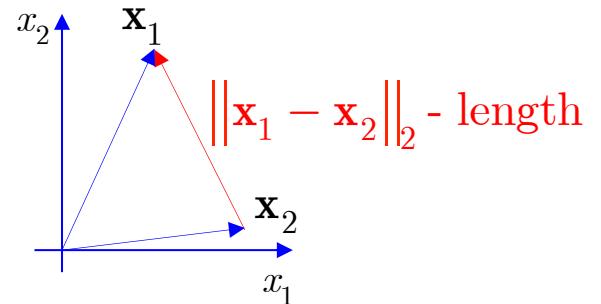
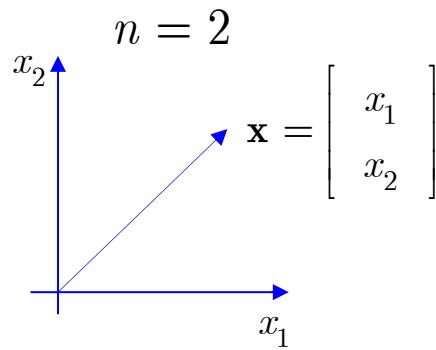
# Distances



$$\begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} N = n \times m$$



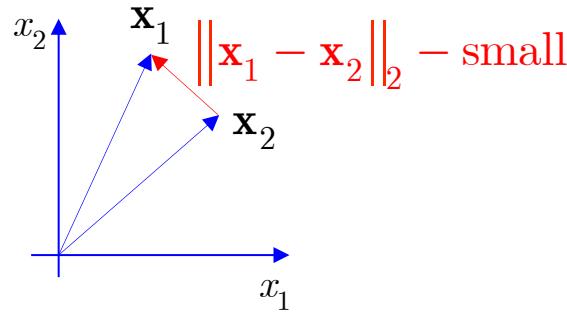
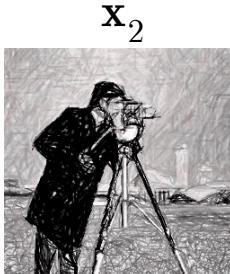
$$d_2(\mathbf{x}_1, \mathbf{x}_2) = \|\mathbf{x}_1 - \mathbf{x}_2\|_2 = \left( \sum_{i=1}^N (x_{1i} - x_{2i})^2 \right)^{\frac{1}{2}}$$



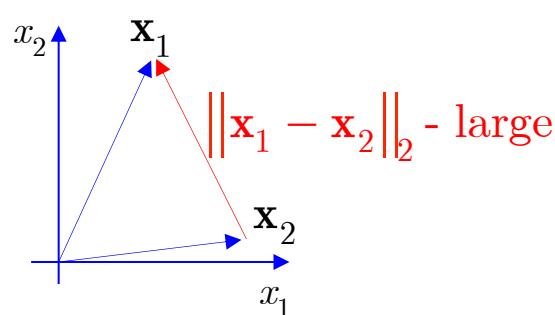
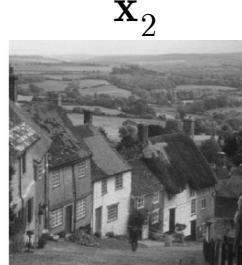
# Distances

---

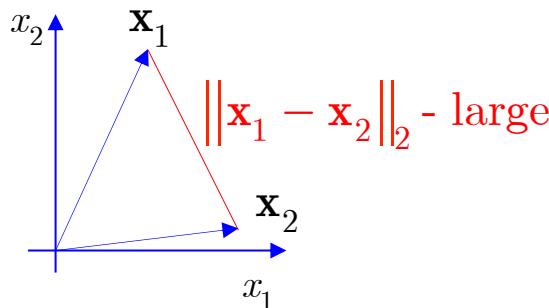
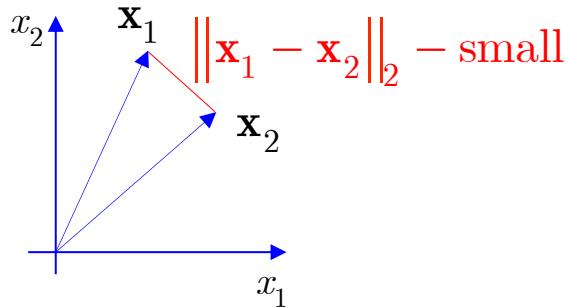
- Perceptually close images



- Perceptually different images



# Distances: link to inner product



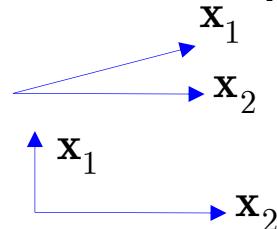
$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \|\mathbf{x}_1\|_2^2 - 2 \underbrace{\mathbf{x}_1^T \mathbf{x}_2}_{\text{inner product}} + \|\mathbf{x}_2\|_2^2$$

Recall B: matrix operations

$$\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} \in \mathbb{R} \triangleq \left[ \begin{array}{cccc} x_1 & x_2 & \dots & x_n \end{array} \right] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \left[ \begin{array}{c} - \\ \mathbf{x} \\ - \end{array} \right] \begin{bmatrix} | \\ \mathbf{y} \\ | \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

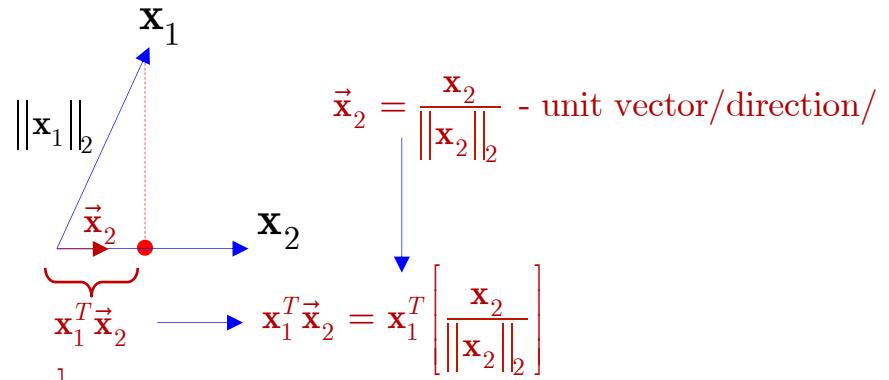
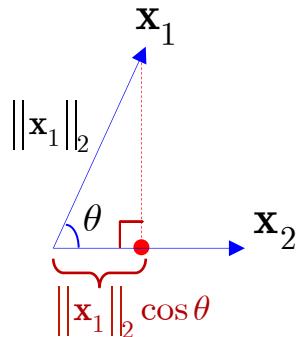
$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2 - \text{small} \Rightarrow \mathbf{x}_1^T \mathbf{x}_2 - \text{large}$$

$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2 - \text{large} \Rightarrow \mathbf{x}_1^T \mathbf{x}_2 - \text{small}$$



# Distances: link to cosine similarity

- Cosine similarity



$$\|x_1\|_2 \cos \theta = x_1^T \vec{x}_2 = x_1^T \left[ \frac{\mathbf{x}_2}{\|\mathbf{x}_2\|_2} \right]$$

$$x_1^T x_2 = \|x_1\|_2 \|x_2\|_2 \cos \theta$$

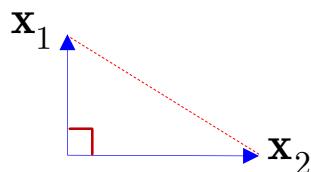
$$\Rightarrow \cos \theta = \frac{x_1^T x_2}{\|x_1\|_2 \|x_2\|_2}$$

# Distances: orthogonal and collinear vectors

---

$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - 2\mathbf{x}_1^T \mathbf{x}_2 = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - 2\cos\theta \|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2$$

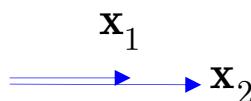
- Orthogonal vectors



$$\begin{aligned}\mathbf{x}_1^T \mathbf{x}_2 &= 0 \\ \|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2 \underbrace{\cos 90^\circ}_{=0} &= 0\end{aligned}$$

$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - 2\underbrace{\mathbf{x}_1^T \mathbf{x}_2}_{=0} = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2$$

- Collinear vectors



$$\mathbf{x}_1^T \mathbf{x}_2 = \max$$

$$\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2 \underbrace{\cos 0^\circ}_{=1} = \|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2$$

$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = \|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 - 2\|\mathbf{x}_1\|_2 \|\mathbf{x}_2\|_2$$

$$\xrightarrow{\mathbf{x}_1 = \mathbf{x}_2}$$

$$\|\mathbf{x}_1\|_2 = \|\mathbf{x}_2\|_2$$

$$\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2 = 0$$

# Transformation of random variables

---

- General linear transform:
  - Random scalar variables
  - Random vectors
- Sum of random variables
  - Sum of independent random variables
  - Concentration inequalities
  - Weak Law of Large Numbers and CLT
  - Application to vector norms (see later after transformation of Gaussian random vectors)
- Transformation of Gaussian random vectors
  - Mean and covariance matrix after transformation
  - Pdf after transformation
  - Properties of transform
    - Decorrelation
    - Whitening

# Second order models

---

- Plan
  - Introduce correlation and covariance matrices
  - Introduce uncorrelated and independent vectors
  - Introduce bivariate and multidimensional Gaussian pdfs
  - Investigate transformation of Gaussian random vectors:
    - pdf and moments

# Second order stochastic models

- **Covariance and cross-covariance** Only for the second order statistics
  - Covariance matrix of vector  $\mathbf{X} \in \mathbb{R}^N$   $\mathbf{K}_{\mathbf{xx}} \in \mathbb{R}^{N \times N}$ 
$$\mathbf{K}_{\mathbf{xx}} = \text{Var}[\mathbf{X}] = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] = E[\mathbf{XX}^T] - E[\mathbf{X}]E[\mathbf{X}]^T$$
  - Cross-covariance matrix between vectors  $\mathbf{X} \in \mathbb{R}^N$  and  $\mathbf{Y} \in \mathbb{R}^M$   $\mathbf{K}_{\mathbf{xy}} \in \mathbb{R}^{N \times M}$ 
$$\mathbf{K}_{\mathbf{xy}} = \text{Cov}[\mathbf{X}, \mathbf{Y}] = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T] = E[\mathbf{XY}^T] - E[\mathbf{X}]E[\mathbf{Y}]^T$$
- Properties
  - The covariance matrix is a **symmetric matrix**  $\mathbf{K}_{\mathbf{xx}}^T = \mathbf{K}_{\mathbf{xx}}$
  - The covariance matrix is a **positive semidefinite matrix**  $\mathbf{a}^T \mathbf{K}_{\mathbf{xx}} \mathbf{a} \geq 0$
  - The transpose of cross-covariance matrix is:  $\mathbf{K}_{\mathbf{xy}}^T = \mathbf{K}_{\mathbf{yx}}$

# Second order stochastic models

---

- Example of covariance matrix

$$\mathbf{K}_{\mathbf{xx}} = \text{Cov}[\mathbf{X}, \mathbf{X}] = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] = E[\mathbf{XX}^T] - E[\mathbf{X}]E[\mathbf{X}]^T =$$
$$\begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \dots & E[(X_1 - E[X_1])(X_N - E[X_N])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \dots & E[(X_2 - E[X_2])(X_N - E[X_N])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_N - E[X_N])(X_1 - E[X_1])] & E[(X_N - E[X_N])(X_2 - E[X_2])] & \dots & E[(X_N - E[X_N])(X_N - E[X_N])] \end{bmatrix}$$

$$\mathbf{K}_{x_i x_j} = \text{Cov}[X_i, X_j] = E[(X_i - E[X_i])(X_j - E[X_j])]$$

# Second order stochastic models

---

- Example of cross-covariance matrix

$$\mathbf{K}_{\mathbf{xy}} = \text{Cov}[\mathbf{X}, \mathbf{Y}] = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{Y} - E[\mathbf{Y}])^T] = E[\mathbf{XY}^T] - E[\mathbf{X}]E[\mathbf{Y}]^T =$$
$$\begin{bmatrix} E[(X_1 - E[X_1])(Y_1 - E[Y_1])] & E[(X_1 - E[X_1])(Y_2 - E[Y_2])] & \dots & E[(X_1 - E[X_1])(Y_M - E[Y_M])] \\ E[(X_2 - E[X_2])(Y_1 - E[Y_1])] & E[(X_2 - E[X_2])(Y_2 - E[Y_2])] & \dots & E[(X_2 - E[X_2])(Y_M - E[Y_M])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_N - E[X_N])(Y_1 - E[Y_1])] & E[(X_N - E[X_N])(Y_2 - E[Y_2])] & \dots & E[(X_N - E[X_N])(Y_M - E[Y_M])] \end{bmatrix}$$

$$\mathbf{K}_{x_i y_j} = Cov[X_i, Y_j] = E[(X_i - E[X_i])(Y_j - E[Y_j])]$$

# Second order stochastic models

---

- **Correlation and cross-correlation** Only for the second order statistics

- Correlation matrix of vector  $\mathbf{X} \in \mathbb{R}^N$

$$\mathbf{R}_{\mathbf{xx}} = \mathbb{E}[\mathbf{XX}^T]$$

$$\mathbf{R}_{\mathbf{xx}} \in \mathbb{R}^{N \times N}$$

- Cross-correlation matrix between vectors  $\mathbf{X} \in \mathbb{R}^N$  and  $\mathbf{Y} \in \mathbb{R}^M$

$$\mathbf{R}_{\mathbf{xy}} = \mathbb{E}[\mathbf{XY}^T]$$

$$\mathbf{R}_{\mathbf{xy}} \in \mathbb{R}^{N \times M}$$

- Properties

- The link between covariance and correlation

$$\mathbf{R}_{\mathbf{xx}} = \mathbf{K}_{\mathbf{xx}} + \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^T$$

- The link between cross-covariance and cross-correlation

$$\mathbf{R}_{\mathbf{xy}} = \mathbf{K}_{\mathbf{xy}} + \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{Y}]^T$$

# Summary: random vectors

---

- **Orthogonal random** vectors  $\mathbf{X} \in \mathbb{R}^N$  and  $\mathbf{Y} \in \mathbb{R}^M$

$$E[\mathbf{XY}^T] = \mathbf{0}_{N \times M}$$

- **Uncorrelated random** vectors  $\mathbf{X} \in \mathbb{R}^N$  and  $\mathbf{Y} \in \mathbb{R}^M$

$$E[\mathbf{XY}^T] = E[\mathbf{X}]E[\mathbf{Y}]^T = \underline{\mu}_X \underline{\mu}_Y^T$$

$$\Rightarrow Cov[\mathbf{X}, \mathbf{Y}] = E[\mathbf{XY}^T] - \underline{\mu}_X \underline{\mu}_Y^T = E[\mathbf{X}]E[\mathbf{Y}]^T - \underline{\mu}_X \underline{\mu}_Y^T = \mathbf{0}$$

Only for the second order statistics

If  $\mathbf{X}$  and  $\mathbf{Y}$  are uncorrelated, then  $\mathbf{X} - \underline{\mu}_X$  and  $\mathbf{Y} - \underline{\mu}_Y$  are orthogonal.

- **Independent random** vectors  $\mathbf{X} \in \mathbb{R}^N$  and  $\mathbf{Y} \in \mathbb{R}^M$

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y}) \Rightarrow E[\mathbf{XY}^T] = E[\mathbf{X}]E[\mathbf{Y}^T] = \underline{\mu}_X \underline{\mu}_Y^T$$

$$E[\mathbf{XY}^T] = \int \int \mathbf{xy}^T p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x} \int \mathbf{y}^T p(\mathbf{y}) d\mathbf{y} = E[\mathbf{X}]E[\mathbf{Y}^T] = \underline{\mu}_X \underline{\mu}_Y^T$$

# Second order models: Bivariate Gaussian

$$f(\mathbf{x}) = p(x_1, x_2) = \frac{1}{2\pi\sigma_{X_1}\sigma_{X_2}\sqrt{1-\rho^2}} \exp \left[ -\frac{1}{2(1-\rho^2)} \left( \frac{(x_1 - \bar{x}_1)^2}{\sigma_{X_1}^2} + \frac{(x_2 - \bar{x}_2)^2}{\sigma_{X_2}^2} - 2\frac{\rho(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)}{\sigma_{X_1}\sigma_{X_2}} \right) \right]$$

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}, \mathbf{K}_{\mathbf{xx}} = \begin{bmatrix} \sigma_{X_1}^2 & \rho\sigma_{X_1}\sigma_{X_2} \\ \rho\sigma_{X_1}\sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix}, \quad \rho = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)\text{Var}(X_2)}}, \quad -1 \leq \rho \leq +1$$

$$\Rightarrow f(\mathbf{x}) = f(x_1, x_2) = \frac{1}{\sqrt{(2\pi)^2 |\det \mathbf{K}_{\mathbf{xx}}|}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right]$$

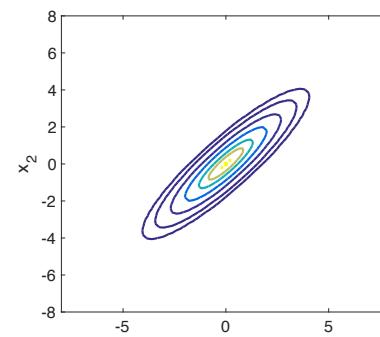
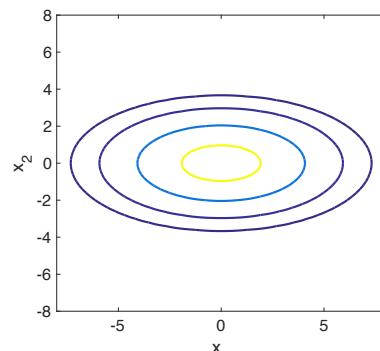
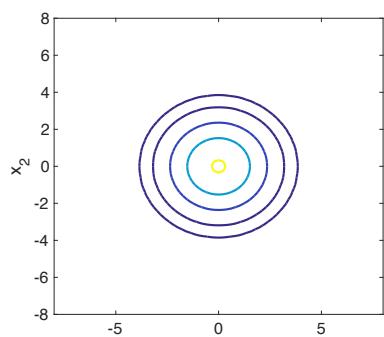
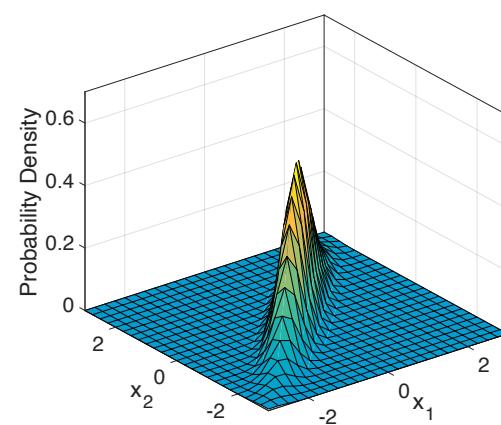
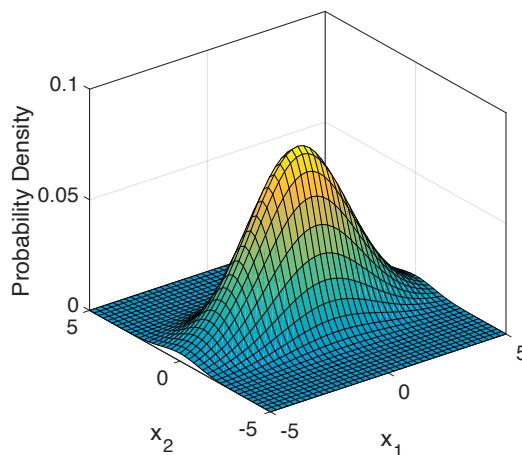
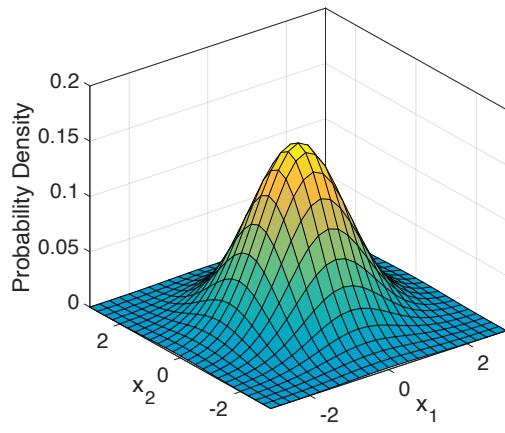
- Uncorrelated data  $\rho = 0$ :  $\mathbf{K}_{\mathbf{xx}} = \begin{bmatrix} \sigma_{X_1}^2 & 0 \\ 0 & \sigma_{X_2}^2 \end{bmatrix}$

$$f(\mathbf{x}) = f(x_1, x_2) = \frac{1}{2\pi\sigma_{X_1}\sigma_{X_2}} \exp \left[ -\frac{1}{2} \left( \frac{(x_1 - \bar{x}_1)^2}{\sigma_{X_1}^2} + \frac{(x_2 - \bar{x}_2)^2}{\sigma_{X_2}^2} \right) \right]$$

$$\Rightarrow f(\mathbf{x}) = f(x_1, x_2) = \underbrace{\frac{1}{\sqrt{2\pi\sigma_{X_1}^2}} \exp \left[ -\frac{(x_1 - \bar{x}_1)^2}{2\sigma_{X_1}^2} \right]}_{f(x_1)} \underbrace{\frac{1}{\sqrt{2\pi\sigma_{X_2}^2}} \exp \left[ -\frac{(x_2 - \bar{x}_2)^2}{2\sigma_{X_2}^2} \right]}_{f(x_2)}$$

It implies independence  $\longrightarrow f(x_1) f(x_2)$

# PCA on synthetic data: bivariate Gaussian



$$\begin{bmatrix} \sigma_{X_1}^2 & 0 \\ 0 & \sigma_{X_2}^2 \end{bmatrix} \sigma_{X_1}^2 = \sigma_{X_2}^2 = 1$$

$$\begin{bmatrix} \sigma_{X_1}^2 & 0 \\ 0 & \sigma_{X_2}^2 \end{bmatrix} \sigma_{X_1}^2 = 4$$

$$\begin{bmatrix} \sigma_{X_1}^2 & \rho\sigma_{X_1}\sigma_{X_2} \\ \rho\sigma_{X_1}\sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix} \begin{array}{l} \sigma_{X_1}^2 = 1 \\ \sigma_{X_2}^2 = 1 \\ \rho = 0.9 \end{array}$$

# Second order models: Multivariate Gaussian

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_N) = \frac{1}{\sqrt{(2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}|}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right]$$

$$(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) = \|(\mathbf{x} - \bar{\mathbf{x}})\|_{\mathbf{K}_{\mathbf{xx}}^{-1}}^2 = D_M \text{ - the Mahalanobis distance}$$

**Notation:**  $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$

**Remark:** we talk about a random vector  $\mathbf{X}$  generated from  $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$

It should not be confused with the data matrix! It should be clear from the text.

For square matrix:  $\mathbf{K}_{\mathbf{xx}} = \mathbf{U}\Sigma\mathbf{U}^{-1}$

- **Eigenvectors**  $\mathbf{U} \in \underbrace{\mathbb{C}^{n \times n}}_{\text{generally complex}} = \begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & & | \end{bmatrix}$
- **Eigenvalues**  $\sigma_1, \dots, \sigma_n, \quad \underbrace{\sigma_i \in \mathbb{C}}_{\text{generally complex}}$   $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$  or  $\Sigma_{ii} = \sigma_i, \forall i$

Recall C: eigenvalue and SVD decompositions

Recall: for square and symmetric  
 $\mathbf{K}_{\mathbf{xx}} = \mathbf{U}\Sigma\mathbf{U}^T$

# Summary of notations

---

$$\mathcal{N}(\mu, \sigma^2)$$

$$X \sim \mathcal{N}(\bar{x}, \sigma_X^2)$$

$$\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{C}_X)$$

$$\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$$

$$\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \sigma_X^2 \mathbf{I}_N)$$

$$\mathcal{L}(\mu, \lambda)$$

$$\mathcal{U}(a, b)$$

$$\mathcal{GGD}(\mu, \gamma, \lambda)$$

$$\eta(\gamma, \lambda)$$

Gaussian pdf with mean  $\mu$  and variance  $\sigma^2$

Scalar random variable with Gaussian pdf with mean  $\bar{x}$  and variance  $\sigma_X^2$

Vector random variable with Gaussian pdf with mean vector  $\bar{\mathbf{x}}$  and covariance matrix  $\mathbf{C}_X/\mathbf{K}_{\mathbf{xx}}$

Vector RV with Gaussian pdf with mean vector  $\bar{\mathbf{x}}$  and diagonal covariance matrix of size  $N$  with  $\sigma_X^2$  on the main diagonal

Laplacian pdf with mean  $\mu$  and scale  $\lambda$

Uniform pdf on interval [a,b]

Generalized Gaussian pdf with mean  $\mu$ , shape parameter  $\gamma$  and scale  $\lambda$

Parameter of  $\mathcal{GGD}$  pdf

# Second order models: Multivariate Gaussian

---

- Generation of data with desired covariance matrix

$$\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$$

$$\mathbf{K}_{\mathbf{xx}} = \mathbf{U}\Sigma\mathbf{U}^T = \mathbf{U}\Sigma^{1/2} \left( \mathbf{U}\Sigma^{1/2} \right)^T$$

$$\begin{aligned} \mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}}) &\Leftrightarrow \mathbf{X} = \bar{\mathbf{x}} + \mathbf{U}\Sigma^{1/2}\mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &\Leftrightarrow \mathbf{X} = \bar{\mathbf{x}} + \mathbf{U}\mathcal{N}(\mathbf{0}, \Sigma) \end{aligned}$$

# Linear transform of multivariate Gaussian vector

---

- Linear transformation:

- Mean
- Covariance matrix
- Given an observation vector  $\mathbf{x}$  generated from  $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$
- Apply a **generic linear transform** to this vector

$$\tilde{\mathbf{x}} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

- **Mean** of transformed vector

$$\tilde{\mathbf{x}} = \mathbf{E}[\tilde{\mathbf{X}}] = \mathbf{E}[\mathbf{W}^T \mathbf{X} + \mathbf{b}] = \mathbf{W}^T \mathbf{E}[\mathbf{X}] + \mathbf{b} = \mathbf{W}^T \bar{\mathbf{x}} + \mathbf{b}$$

Remark: The Gaussian random vector under the linear (affine) transform will remain Gaussian but with modified mean and covariance.

Therefore, it suffices to find new mean and covariance matrix.

# Linear transform of multivariate Gaussian vector

- Linear transformation:

- Mean
- Covariance matrix
- **Covariance matrix** of transformed vector

$$\begin{aligned} \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} &= \mathbf{E}\left[\left(\tilde{\mathbf{X}} - \mathbf{E}(\tilde{\mathbf{X}})\right)\left(\tilde{\mathbf{X}} - \mathbf{E}(\tilde{\mathbf{X}})\right)^T\right] \\ &\quad \tilde{\mathbf{X}} - \mathbf{E}(\tilde{\mathbf{X}}) = \mathbf{W}^T \mathbf{X} + \mathbf{b} - \mathbf{W}^T \bar{\mathbf{x}} - \mathbf{b} = \mathbf{W}^T (\mathbf{X} - \bar{\mathbf{x}}) \\ \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} &= \mathbf{E}\left[\left(\mathbf{W}^T (\mathbf{X} - \bar{\mathbf{x}})\right)\left(\mathbf{W}^T (\mathbf{X} - \bar{\mathbf{x}})\right)^T\right] = \mathbf{E}\left[\mathbf{W}^T (\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T \mathbf{W}\right] \\ &= \mathbf{W}^T \underbrace{\mathbf{E}\left[(\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T\right]}_{\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \mathbf{U}\Sigma\mathbf{U}^T} \mathbf{W} = \mathbf{W}^T \mathbf{C} = \mathbf{W}^T \mathbf{U} \Sigma \mathbf{U}^T \mathbf{W} \end{aligned}$$

- If  $\mathbf{W}^T = \mathbf{U}^T$ , i.e., the transform is chosen to be the **decorrelation transform**

$$\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \mathbf{W}^T \mathbf{U} \Sigma \mathbf{U}^T \mathbf{W} = \underbrace{\mathbf{U}^T}_{\mathbf{I}} \mathbf{U} \Sigma \underbrace{\mathbf{U}^T}_{\mathbf{I}} \mathbf{U} = \Sigma = \text{diag}\{\sigma_1, \dots, \sigma_n\}$$

# Linear transform of multivariate Gaussian vector

- **Decorrelation transform:**

- Given an observation vector  $\mathbf{x}$  generated from  $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$
- Apply the transform to the centered vector:

$$\tilde{\mathbf{x}} = \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$$

with  $\mathbf{U}$  to be a eigenvector matrix of  $\mathbf{K}_{\mathbf{xx}} = \mathbf{U}\Sigma\mathbf{U}^T$

- The transformed data are also Gaussian  $\tilde{\mathbf{X}} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}})$  with the covariance matrix:

$$\begin{aligned} \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} &= \mathbf{E}[\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T] = \mathbf{E}\left[\mathbf{U}^T(\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T \mathbf{U}\right] = \mathbf{U}^T \underbrace{\mathbf{E}\left[(\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T\right]}_{\mathbf{K}_{\mathbf{xx}}} \mathbf{U} \\ &= \mathbf{U}^T \mathbf{K}_{\mathbf{xx}} \mathbf{U} = \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} \underbrace{\Sigma}_{\mathbf{I}} \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} = \Sigma \end{aligned}$$

$$\mathbf{K}_{\mathbf{xx}} = \begin{bmatrix} \sigma_{X_1}^2 & \boxed{\rho\sigma_{X_1}\sigma_{X_2}} \\ \boxed{\rho\sigma_{X_1}\sigma_{X_2}} & \sigma_{X_2}^2 \end{bmatrix} \xrightarrow{\mathbf{U}^T} \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

decorrelation

$$\rho \text{ any } \xrightarrow{\mathbf{U}^T} \tilde{\rho} = 0$$

# Linear transform of multivariate Gaussian vector

- **Energy compaction (EC) property:**

$$\mathbf{K}_{\mathbf{x}\mathbf{x}} = \begin{bmatrix} \sigma_{X_1}^2 & \rho\sigma_{X_1}\sigma_{X_2} \\ \rho\sigma_{X_1}\sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix} \xrightarrow{\mathbf{U}^T} \mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$$

All “energy” or information is only in the diagonal elements

$$\sigma_{X_1}^2, \sigma_{X_2}^2 \text{ any } \xrightarrow{\mathbf{U}^T} \sigma_1 \geq \sigma_2$$

Decaying eigenvalues

# Linear transform of multivariate Gaussian vector

---

- **Whitening property :**

- Given an observation vector  $\mathbf{x}$  generated from  $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$
- Apply **whitening** PCA to the centered vector:

$$\tilde{\mathbf{x}} = \Sigma^{-1/2} \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$$

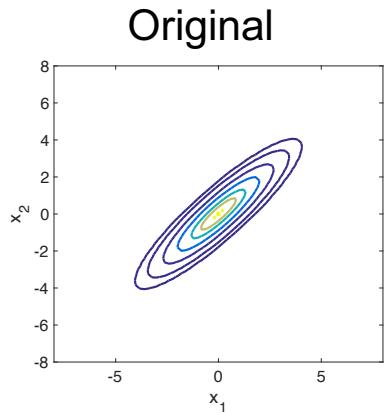
with  $\Sigma^{-1/2} = \text{diag}\left(1/\sqrt{\sigma_i}\right)$

- The transformed data are also Gaussian with the covariance matrix:  $\tilde{\mathbf{X}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\begin{aligned}\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} &= \mathbf{E}[\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T] = \mathbf{E}\left[\Sigma^{-1/2}\mathbf{U}^T(\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T\mathbf{U}\Sigma^{-1/2}\right] \\ &= \Sigma^{-1/2}\mathbf{U}^T \underbrace{\mathbf{E}\left[(\mathbf{X} - \bar{\mathbf{x}})(\mathbf{X} - \bar{\mathbf{x}})^T\right]}_{\mathbf{K}_{\mathbf{xx}}} \mathbf{U}\Sigma^{-1/2} \\ &= \Sigma^{-1/2}\mathbf{U}^T \mathbf{K}_{\mathbf{xx}} \mathbf{U}\Sigma^{-1/2} = \underbrace{\Sigma^{-1/2}\mathbf{U}^T}_{\mathbf{I}} \underbrace{\mathbf{U}\Sigma}_{\mathbf{I}} \underbrace{\mathbf{U}^T}_{\mathbf{I}} \underbrace{\mathbf{U}\Sigma^{-1/2}}_{\mathbf{I}} = \Sigma^{-1/2}\Sigma\Sigma^{-1/2} = \mathbf{I}\end{aligned}$$

# Demo

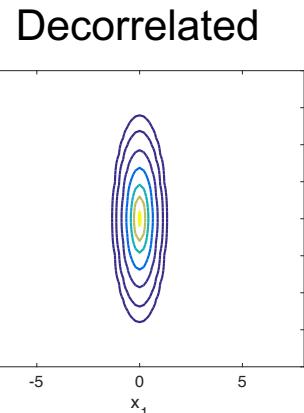
- Matlab demo  $N = 2$



$\mathbf{x}$

$$\mathbf{K}_{\mathbf{x}\mathbf{x}} = \mathbf{U}\Sigma\mathbf{U}^T$$

$$\mathbf{K}_{\mathbf{x}\mathbf{x}} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}, \rho = 0.9$$

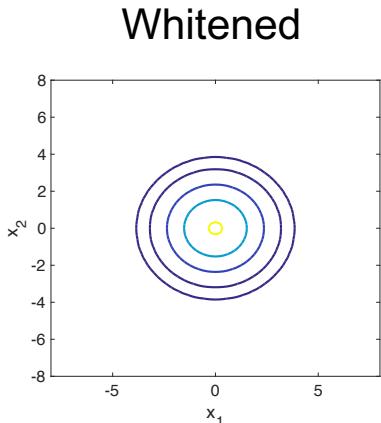


$$\tilde{\mathbf{x}} = \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$$

$$\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \Sigma$$

$$\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \begin{pmatrix} 0.1 & 0 \\ 0 & 1.9 \end{pmatrix}$$

Note: Matlab ordering

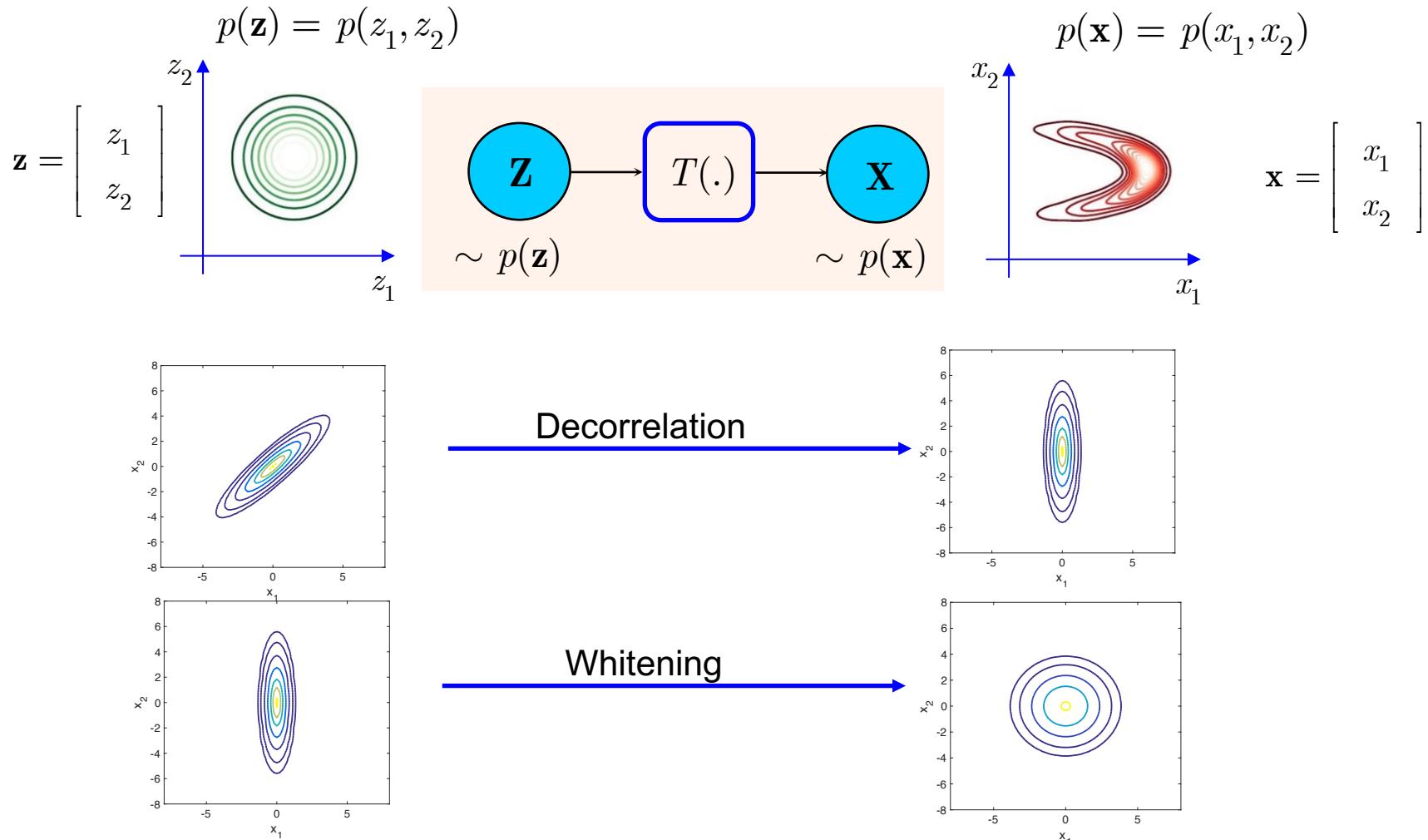


$$\tilde{\mathbf{x}} = \Sigma^{-1/2} \mathbf{U}^T (\mathbf{x} - \bar{\mathbf{x}})$$

$$\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \mathbf{I}$$

$$\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

# Recall: overall goal behind transformation



# N-order Markov sequences

- The  $N$ -order stationary Markov sequence:

$$\mathbf{K}_{\mathbf{xx}} = \sigma_X^2 \begin{bmatrix} 1 & \rho & \rho^2 \dots & \rho^{N-1} \\ \rho & 1 & \dots & \rho^{N-2} \\ \vdots & \vdots & \ddots & \vdots \\ \rho^{N-1} & \rho^{N-2} & \dots & 1 \end{bmatrix} \begin{aligned} \mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}}) &\Leftrightarrow \mathbf{X} = \bar{\mathbf{x}} + \mathbf{U}\Sigma^{1/2}\mathcal{N}(\mathbf{0}, \mathbf{I}) \\ &\Leftrightarrow \mathbf{X} = \bar{\mathbf{x}} + \mathbf{U}\mathcal{N}(\mathbf{0}, \Sigma) \end{aligned}$$

$$N = 4 \quad \rho = 0.95$$

$$\mathbf{K}_{\mathbf{xx}} = \sigma_X^2 \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

$\mathbf{U} =$  High frequency      Low frequency

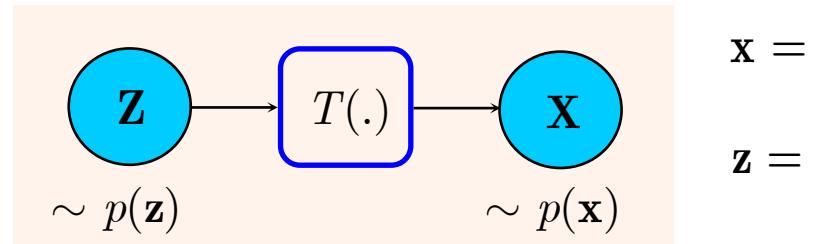
|         |         |         |        |
|---------|---------|---------|--------|
| 0.2747  | -0.5062 | 0.6516  | 0.4937 |
| -0.6516 | 0.4937  | 0.2747  | 0.5062 |
| 0.6516  | 0.4937  | -0.2747 | 0.4937 |
| -0.2747 | -0.5062 | -0.6516 | 0.4937 |

$\mathbf{D} =$

|        |        |        |        |
|--------|--------|--------|--------|
| 0.0300 | 0      | 0      | 0      |
| 0      | 0.0506 | 0      | 0      |
| 0      | 0      | 0.1627 | 0      |
| 0      | 0      | 0      | 3.7568 |

# Linear transform: a link to affine transform

---



$$\begin{aligned}\mathbf{x} &= \underbrace{\mathbf{W}^T}_{\mathbf{A}} \mathbf{z} = \mathbf{A}\mathbf{z} \\ \mathbf{z} &= \mathbf{A}^{-1}\mathbf{x}\end{aligned}$$

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(T^{-1}(\mathbf{x}))}{|\det J_T(\mathbf{z})|} = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det J_T(\mathbf{z})|} = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det \mathbf{A}|}$$

- Consider  $p(\mathbf{z}) = \mathcal{N}(\bar{\mathbf{z}}, \mathbf{K}_{\mathbf{zz}})$
- Find  $p(\mathbf{x})$ ?

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det \mathbf{A}|} = \frac{1}{|\det \mathbf{A}|} \frac{1}{\sqrt{(2\pi)^N |\det \mathbf{K}_{\mathbf{zz}}|}} \exp \left[ -\frac{1}{2} (\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}})^T \mathbf{K}_{\mathbf{zz}}^{-1} (\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}}) \right]$$

# Linear transform: a link to affine transform

---

- Simplification

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(\mathbf{A}^{-1}\mathbf{x})}{|\det \mathbf{A}|} = \frac{1}{|\det \mathbf{A}|} \frac{1}{\sqrt{(2\pi)^N |\det \mathbf{K}_{zz}|}} \exp \left[ -\frac{1}{2} (\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}})^T \mathbf{K}_{zz}^{-1} (\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}}) \right]$$

$$(\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}}) = \mathbf{A}^{-1}(\mathbf{x} - \mathbf{A}\bar{\mathbf{z}})$$

$$(\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}})^T = (\mathbf{A}^{-1}(\mathbf{x} - \mathbf{A}\bar{\mathbf{z}}))^T = (\mathbf{x} - \mathbf{A}\bar{\mathbf{z}})^T \mathbf{A}^{-1T}$$

$$(\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}})^T \mathbf{K}_{zz}^{-1} (\mathbf{A}^{-1}\mathbf{x} - \bar{\mathbf{z}}) = (\mathbf{x} - \mathbf{A}\bar{\mathbf{z}})^T (\mathbf{A}\mathbf{K}_{zz}\mathbf{A}^T)^{-1} (\mathbf{x} - \mathbf{A}\bar{\mathbf{z}})$$

$$\mathbf{A}^{-1T} \mathbf{K}_{zz}^{-1} \mathbf{A}^{-1} = (\mathbf{A}\mathbf{K}_{zz}\mathbf{A}^T)^{-1} = \mathbf{K}_{xx}^{-1}$$

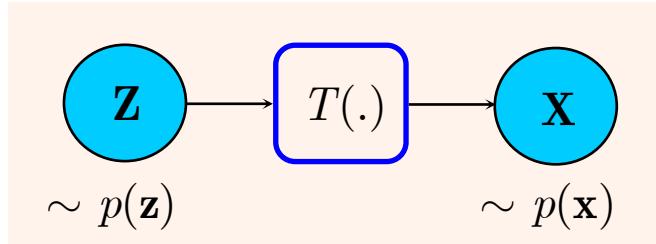
$$|\det \mathbf{K}_{xx}| = |\det \mathbf{A}\mathbf{K}_{zz}\mathbf{A}^T| = |\det \mathbf{A}| |\det \mathbf{K}_{zz}| |\det \mathbf{A}^T| = |\det \mathbf{A}|^2 |\det \mathbf{K}_{zz}|$$

$$\bar{\mathbf{x}} = \mathbf{A}\bar{\mathbf{z}}$$

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} |\det(\mathbf{K}_{xx})|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{xx}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right]$$

# Beyond second order models

---



- We considered generalized linear transform that targets the second order statistics

$$\mathbf{x} = \mathbf{W}^T \mathbf{z} + \mathbf{b}$$

- Possible extension is a non-linear transform

$$\mathbf{x} = \sigma\left(\mathbf{W}^T \mathbf{z} + \mathbf{b}\right)$$

- Nested transforms = deep models

$$\mathbf{x} = \sigma_L\left(\mathbf{W}_L^T \dots \sigma_1\left(\mathbf{W}_1^T \mathbf{z} + \mathbf{b}_1\right) + \mathbf{b}_L\right)$$

$$\Leftrightarrow \mathbf{x} = g_{\theta}(\mathbf{z}) \text{ with } \theta = (\mathbf{W}_1^T, \dots, \mathbf{W}_L^T; \mathbf{b}_1, \dots, \mathbf{b}_L)$$

# Concentration properties and Gaussianity

---

- Concentration plays an important role for the analysis and design of systems
- Plan
  - Sum of jointly Gaussian random variables
    - As a special case of linear transform
  - Sum of squared random variables
    - The Euclidean norm of vector and its concentration
      - Mean
      - Variance
  - Hamming weight as a norm of binary sequences
    - Mean
    - Variance
  - Note on visualization

# Linear transform: a link to affine transform

---

- A special case of linear transformation is

$$X = a_1 Z_1 + a_2 Z_2 + \cdots + a_N Z_N$$

where  $Z_1, Z_2, \dots, Z_N$  are jointly Gaussian

- $X$  can be written

$$X = \begin{bmatrix} a_1 & \cdots & a_N \end{bmatrix} \begin{bmatrix} Z_1 \\ \vdots \\ Z_N \end{bmatrix} \triangleq \mathbf{A}\mathbf{Z}$$

- $X$  is a linear transformation of Gaussians that is also Gaussian

# Linear transform: a link to affine transform

---

- Therefore,  $X$  is Gaussian with **mean**

$$E[X] = E[a_1Z_1 + a_2Z_2 + \cdots + a_NZ_N] = \sum_{i=1}^N a_i \underbrace{E[Z_i]}_{\mu_Z} = \mu_Z \sum_{i=1}^N a_i$$

where  $Z_1, Z_2, \dots, Z_N$  are jointly Gaussian

- Variance:**

$$\text{Var}[X] = \sum_{i=1}^N \sum_{j=1}^N a_i a_j \text{Cov}[Z_i, Z_j]$$

- If  $Z_1, Z_2, \dots, Z_N$  are independent Gaussian variables, i.e.,  $\text{Cov}[Z_i, Z_j] = 0$
- Then the variance of  $X$  reduces to

$$\text{Var}[X] = \sum_{i=1}^N a_i^2 \text{Var}[Z_i]$$

# Concentration properties and Gaussianity

---

- Concentration plays an important role for the analysis and design of systems
- Plan
  - Sum of jointly Gaussian random variables
    - As a special case of linear transform
  - Sum of squared random variables
    - The Euclidean norm of vector and its concentration
      - Mean
      - Variance
  - Hamming weight as a norm of binary sequences
    - Mean
    - Variance
  - Note on visualization

# Concentration: $\ell_2$ - norm

- **Concentration of norm around the mean:**

- Given an i.i.d. **unit variance** random Gaussian vector  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- The Euclidean norm moments

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad S_N(X) = \left\| \mathbf{X} \right\|_2^2 = X_1^2 + X_2^2 + \cdots + X_N^2 \sim \chi^2 \text{ with } N \text{ degrees of freedom}$$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma_X^2 \mathbf{I}) \quad S_N(X) = \left\| \mathbf{X} \right\|_2^2 = X_1^2 + X_2^2 + \cdots + X_N^2 \sim \text{Gamma} \text{ with } N \text{ degrees of freedom}$$

- **Mean:**

$$E\left[\left\| \mathbf{X} \right\|_2^2\right] = E\left[X_1^2 + \cdots + X_N^2\right] = \sum_{i=1}^N \underbrace{E\left[X_i^2\right]}_{\sigma_X^2 = 1} = N$$
$$N\sigma_X^2$$

- **Variance:**

$$\text{Var}\left[\left\| \mathbf{X} \right\|_2^2\right] = E\left[\left(\left\| \mathbf{X} \right\|_2^2\right)^2\right] - E\left[\left\| \mathbf{X} \right\|_2^2\right]^2$$

# Concentration: $\ell_2$ - norm

$$E\left[\left(\|\mathbf{X}\|_2^2\right)^2\right] = E\left[\left(\sum_{i=1}^N X_i^2\right)^2\right] = E\left[\sum_{i=1}^N \sum_{j=1}^N X_i^2 X_j^2\right] = \sum_{i=1}^N \sum_{j=1}^N E[X_i^2 X_j^2]$$

$$\begin{aligned} &= \underbrace{\sum_{i=1}^N E[X_i^4]}_{3N(\sigma_X^2)^2} + 2 \underbrace{\sum_{i=1}^{N-1} \sum_{j=i}^N E[X_i^2] E[X_j^2]}_{\substack{(N-1)\frac{N}{2} \\ \text{sum of triangle matrix elements}}} = 3N + N(N-1) = N(N+2) \\ &\quad 2 \quad (\sigma_{X_i}^2 \sigma_{X_j}^2 + \sigma_{X_{ij}}^2)_{\sigma_X^2=1}=0 \end{aligned}$$

$$3N(\sigma_X^2)^2 + N(N-1)(\sigma_X^2)^2 = N(N+2)(\sigma_X^2)^2$$

$$Var\left\|\|\mathbf{X}\|_2^2\right\| = E\left[\left(\|\mathbf{X}\|_2^2\right)^2\right] - E\left[\left(\|\mathbf{X}\|_2^2\right)^2\right]^2 = N(N+2) - (N)^2 = 2N \\ 2N\sigma_X^2$$

$$\frac{\text{Mean}}{\sqrt{\text{Variance}}} = \frac{E\left[\|\mathbf{X}\|_2^2\right]}{\sqrt{Var\left\|\|\mathbf{X}\|_2^2\right\|}} = \frac{N}{\sqrt{2N}} = \sqrt{\frac{N}{2}}$$

Mean growths faster than std

# Summary: Euclidian norm concentration

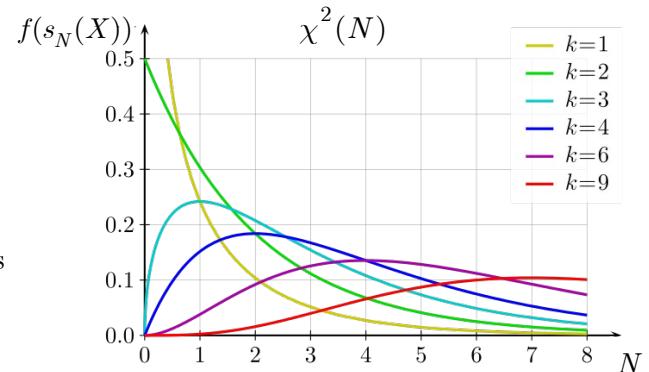
$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{x} = (x_1, \dots, x_N) \quad \|\mathbf{x}\|_2 =$$

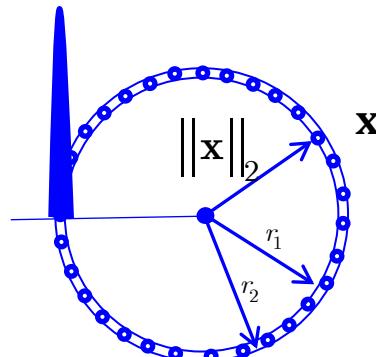
$$\sqrt{\sum_{i=1}^N x_i^2}$$

$x_1^2 + x_2^2 + \dots + x_N^2$

sum of  $N \chi^2(N)$  unit variance R.V.s  
 $\Rightarrow$  CLT  $\Rightarrow$  Gaussian



$$\left\{ \mathbf{x} : r_1 \leq \|\mathbf{x}\|_2 \leq r_2 \right\}$$



$$\begin{array}{c} \rightarrow \leftarrow \\ E[\|\mathbf{x}\|_2^2] \quad Var[\|\mathbf{x}\|_2^2] \end{array}$$

**WLLN in action**

$$\frac{\text{Mean}}{\sqrt{\text{Variance}}} = \frac{E[\|\mathbf{x}\|_2^2]}{\sqrt{Var[\|\mathbf{x}\|_2^2]}} = \frac{N}{\sqrt{2N}} = \sqrt{\frac{N}{2}}$$

# Concentration: $\ell_2$ - norm

- Concentration of norm around the mean:

- i.i.d.  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- Squared sum:

$$S_N(X) = \|\mathbf{X}\|_2^2 \sim \chi^2(N)$$

- Mean:

$$\text{Mean} = E\left[\|\mathbf{X}\|_2^2\right] = N$$

- Variance:

$$\text{Var}\left[\|\mathbf{X}\|_2^2\right] = 2N$$

- i.i.d.  $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \sigma_X^2 \mathbf{I})$

- Squared sum:

$$S_N(X) = \|\mathbf{X}\|_2^2 \sim \text{Gamma}(K, \theta) \\ K = N/2, \theta = 2\sigma_X^2$$

- Mean:

$$\text{Mean} = E\left[\|\mathbf{X}\|_2^2\right] = K\theta = \frac{N}{2}2\sigma_X^2 = N\sigma_X^2$$

- Variance:

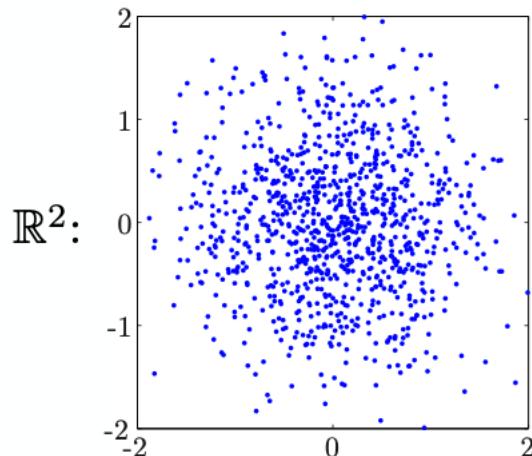
$$\text{Var}\left[\|\mathbf{X}\|_2^2\right] = K\theta^2 = \frac{N}{2}(2\sigma_X^2)^2 = 2N(\sigma_X^2)^2$$

$$\frac{\text{Mean}}{\sqrt{\text{Variance}}} = \frac{E\left[\|\mathbf{X}\|_2^2\right]}{\sqrt{\text{Var}\left[\|\mathbf{X}\|_2^2\right]}} = \frac{N}{\sqrt{2N}} = \sqrt{\frac{N}{2}}$$

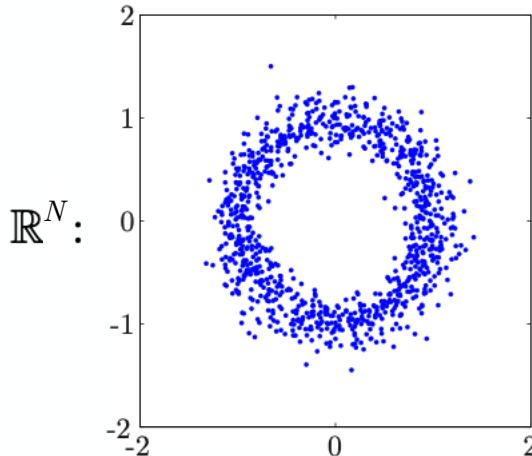
$$\frac{\text{Mean}}{\sqrt{\text{Variance}}} = \frac{E\left[\|\mathbf{X}\|_2^2\right]}{\sqrt{\text{Var}\left[\|\mathbf{X}\|_2^2\right]}} = \frac{N\sigma_X^2}{\sqrt{2N(\sigma_X^2)^2}} = \sqrt{\frac{N}{2}}$$

# Concentration: $\ell_2$ - norm

- To visualize the effect of concentration in high dimensions assume  $\mathbf{X} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{N} \mathbf{I}\right)$



$\mathbb{R}^2:$



Here, only 2 first dim are visualized

- Vectors concentrate near the unit sphere  $E\left[\|\mathbf{X}\|_2^2\right] = N \text{Var}[X_i] = N \frac{1}{N} = 1$
- Here, the variance decreases with the increase of dimension
  - Note: if the dimension increases, our model assumes the decrease of variance inversely proportionally to the dimension. Thus, with  $N \rightarrow \infty$ , we get random Gaussian vectors to be almost “deterministic” and near the unit sphere.

# Concentration properties and Gaussianity

---

- Concentration plays an important role for the analysis and design of systems
- Plan
  - Sum of jointly Gaussian random variables
    - As a special case of linear transform
  - Sum of squared random variables
    - The Euclidean norm of vector and its concentration
      - Mean
      - Variance
  - Hamming weight as a norm of binary sequences
    - Mean
    - Variance
  - Note on visualization

# Recall: Binary or Bernoulli random variable

- Binary random variable takes value from the binary alphabet  $X \in \mathcal{X} = \{0, 1\}$

$$\theta \in [0, 1] \quad p_X(x) = \begin{cases} 1 - \theta, & \text{if } x = 0, \\ \theta, & \text{if } x = 1, \\ 0, & \text{otherwise.} \end{cases} \Rightarrow \Pr[X_i = 1] = \theta, \Pr[X_i = 0] = 1 - \theta$$
$$\Rightarrow \Pr[X = x] = \theta^x (1 - \theta)^{1-x}$$

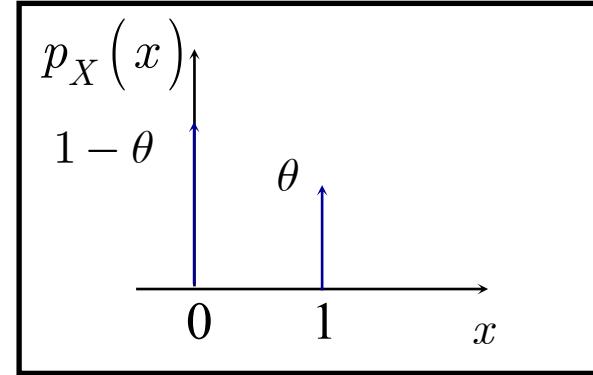
- **Moments**

$$E[X] = \theta \quad \text{Var}[X] = \theta(1 - \theta)$$

- **Proof**

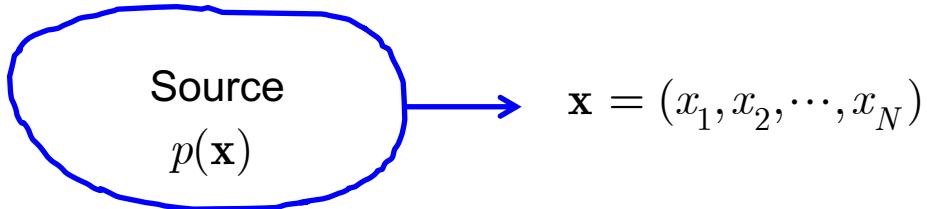
$$E[X] = 0p_X(0) + 1p_X(1) = 0(1 - \theta) + 1\theta = \theta$$

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 = (0^2(1 - \theta) + 1^2\theta) - \theta^2 = \\ &= \theta - \theta^2 = \theta(1 - \theta) \end{aligned}$$



# Binomial pmf: as a sum of iid binary RVs

- Consider a source of iid binary random variables



A **binary memoryless source (BMS)** produces a binary sequence  $\mathbf{x} = (x_1, x_2, \dots, x_N)$  with probability:

$$\Pr[(X_1, X_2, \dots, X_N) = (x_1, x_2, \dots, x_N)] = \prod_{i=1}^N \Pr[X_i = x_i]$$

with probabilities  $\Pr[X_i = 1] = \theta$ ,  $\Pr[X_i = 0] = 1 - \theta$  for  $0 \leq \theta < 1$  and  $1 \leq i \leq N$ .

# Binomial pmf: as a sum of iid binary RVs

- Question: what is the probability to observe a sequence  $\mathbf{x} = (x_1, x_2, \dots, x_N)$ ?

- Given:

BMS

$$\mathbf{x} = (x_1, x_2, \dots, x_N)$$

$$\Pr[X_i = 1] = \theta, \Pr[X_i = 0] = 1 - \theta$$

- Probability that the BMS generates a **specific** sequence  $(x_1, x_2, \dots, x_N)$  that contains ones (Hamming weight  $w$ ) (and  $(N - w)$  zeros) is:

$$\theta^w (1 - \theta)^{N-w}$$

- Probability that  $w$  occurs in a sequence  $(x_1, x_2, \dots, x_N)$  in **any order** is:

$$\Pr[W = w] = \binom{N}{w} \theta^w (1 - \theta)^{N-w}$$

in total there are  $\binom{N}{w} = \frac{N!}{(N-w)!w!}$  sequences containing  $w$  ones.

# Binomial pmf: as a sum of iid binary RVs

- Binomial pmf reflects the statistical behavior of sum of iid binary random variables
- It is also equivalent to the “norm” of binary vector = Hamming weight

$$\mathbf{X} = (X_1, X_2, \dots, X_N) = (0, \underbrace{1}_{\text{with prob } \theta}, \dots, \underbrace{1}_{\text{with prob } \theta})$$

**Hamming weight**  $W = \sum_{i=1}^N X_i \quad W \sim p_W(w) = \mathcal{B}(N, \theta)$

- **Moments**

$$E[W] = E\left[\sum_{i=1}^N X_i\right] = \sum_{i=1}^N \underbrace{E[X_i]}_{\theta} = N\theta$$

$$Var[W] = \sum_{i=1}^N \underbrace{Var[X_i]}_{=\theta(1-\theta)} = N\theta(1-\theta)$$

# Binomial pmf: as a sum of iid binary RVs

- Assume  $\theta = 1/2$  and  $\theta = 1/4$   
for  $N = 32, 1024, 5024$
- Compute probability  $\Pr[W = w]$  for  $w = 0, 1, \dots, N$  as binomial pmf

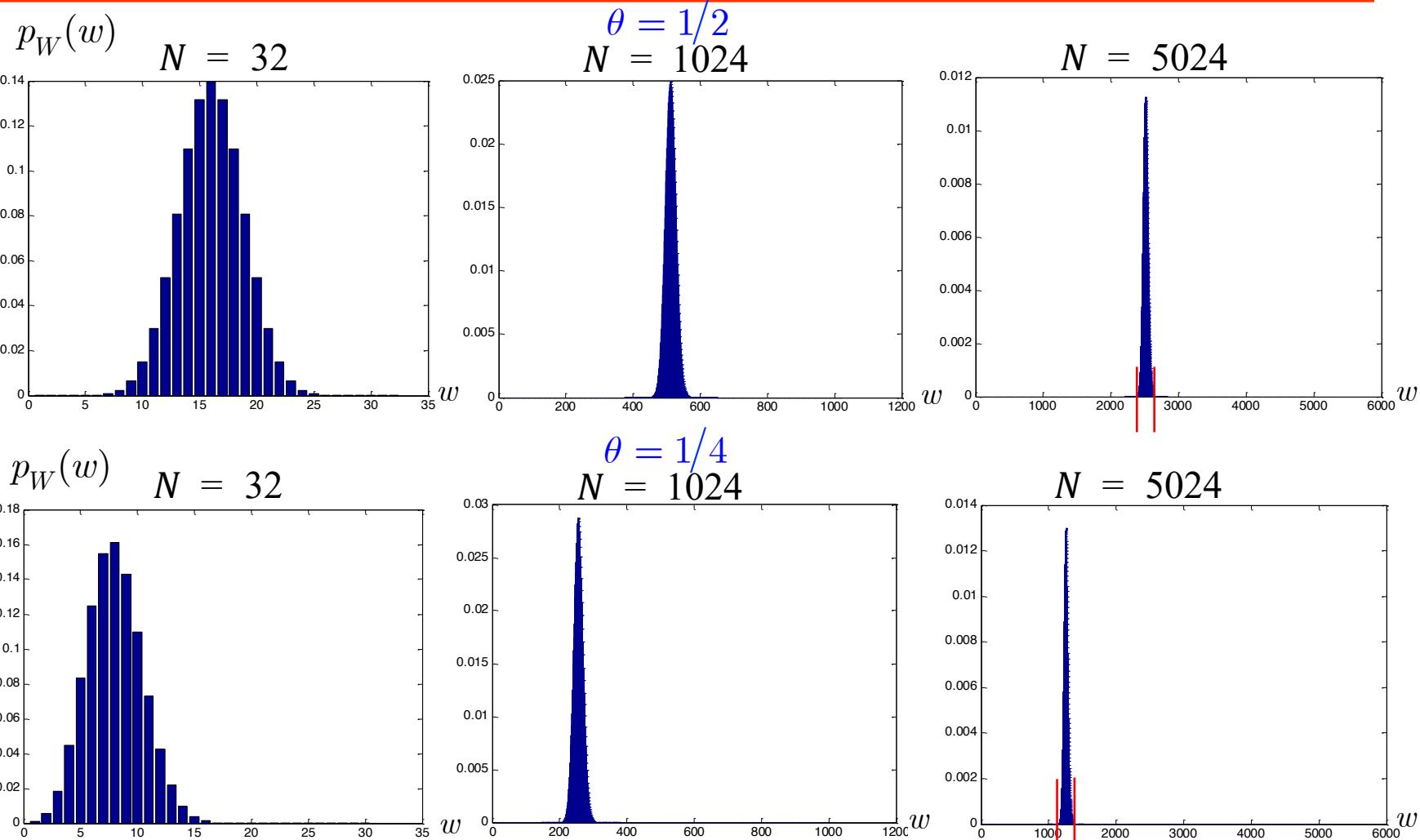
$$\Pr[W = w] = \underbrace{\binom{N}{w}}_{\text{\# of } N\text{-length sequences with } w \text{ ones}} \times \underbrace{\theta^w (1 - \theta)^{N-w}}_{\text{probability to observe } w \text{ ones in the sequence of length } N}$$

**WLLN in action**

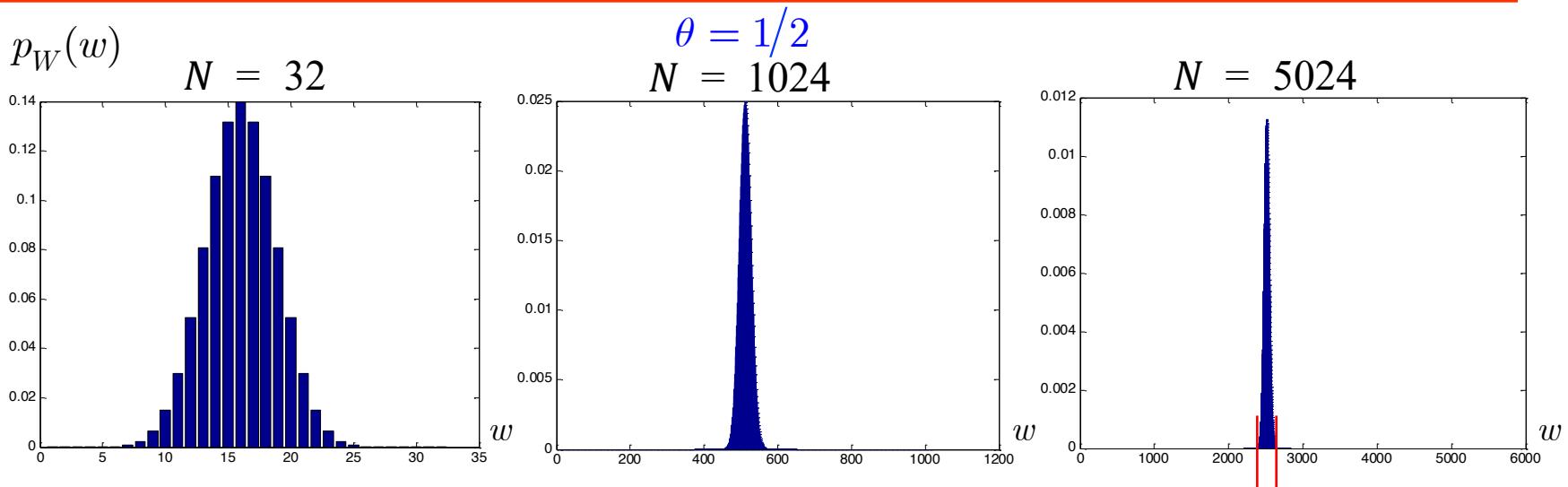
$$\frac{\text{Mean}}{\sqrt{\text{Variance}}} = \frac{E[W]}{\sqrt{Var[W]}} = \frac{N\theta}{\sqrt{N\theta(1 - \theta)}} = \sqrt{\frac{N\theta}{(1 - \theta)}}$$

**Remarks:** pay attention to mean, variance and probability concentration

# Binomial pmf: as a sum of iid binary RVs



# Binomial pmf: as a sum of iid binary RVs



## WLLN in action

- **Remarks:** if the sequence length  $N$  increases:
  - the distribution of  $W$  peaks around  $\theta N = 0.5N$  (mean  $\Rightarrow$  WLLN)
  - the variance of distribution is  $N\theta(1 - \theta)$   $\Rightarrow$  pmf is peaky!

# Concentration properties and Gaussianity

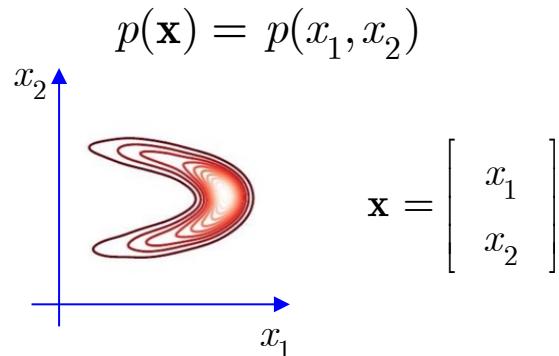
---

- Concentration plays an important role for the analysis and design of systems
- Plan
  - Sum of jointly Gaussian random variables
    - As a special case of linear transform
  - Sum of squared random variables
    - The Euclidean norm of vector and its concentration
      - Mean
      - Variance
  - Hamming weight as a norm of binary sequences
    - Mean
    - Variance
  - Note on visualization

# Visualization of high dimensional data

---

- We can visualize data in 2D or 3D spaces



- Visualization of multidimensional data
  - Reduce dimensions to 2D by special transforms containing the most important data components

$$\tilde{\mathbf{x}} = \mathbf{W}^T \mathbf{x}$$
$$\tilde{\mathbf{x}} \in \mathbb{R}^2 \xleftarrow{\hspace{1cm}} \mathbf{x} \in \mathbb{R}^N$$

- Principle component analysis (PCA)
- T-sne

# Visualization of high dimensional data: T-sne

Journal of Machine Learning Research 9 (2008) 2579-2605

Submitted 5/08; Revised 9/08; Published 11/08

## Stochastic Neighbor Embedding

### Visualizing Data using t-SNE

**Laurens van der Maaten**

TiCC

Tilburg University

P.O. Box 90153, 5000 LE Tilburg, The Netherlands

LVDMAATEN@GMAIL.COM

**Geoffrey Hinton**

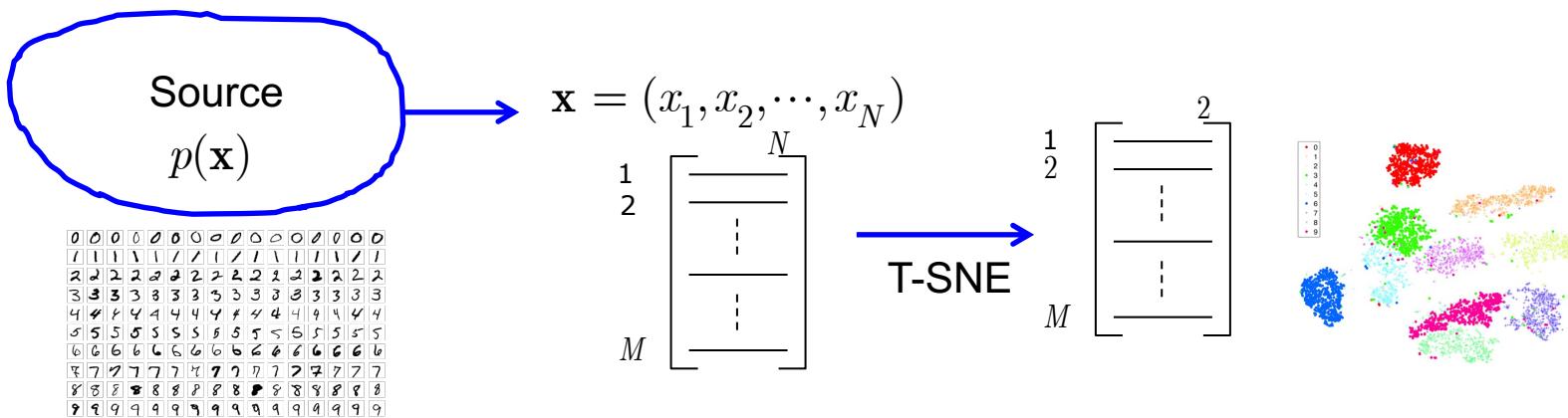
Department of Computer Science

University of Toronto

6 King's College Road, M5S 3G4 Toronto, ON, Canada

HINTON@CS.TORONTO.EDU

<http://jmlr.csail.mit.edu/papers/volume9/vandermaaten08a/vandermaaten08a.pdf>



# Appendix

---

- ▶ Recall A: Basic variance properties

# Recall of variance and its properties

**Definition (Variance of random variable):**

$$Var[X] = \sigma_X^2 = E_{p(x)}[(X - \mu_X)^2]$$

where

$$\mu_X = E_{p(x)}[X] \quad Var[X] = \sum_{x \in S_X} (x - \mu_X)^2 p_X(x)$$

- Standard deviation (Fr: écart-type)

$$\sigma_X = \sqrt{\sigma_X^2}$$

# Recall of variance and its properties

- Two forms of variance definition

$$\text{Var}[X] = E[X^2] - (E[X])^2 = E[X^2] - \mu_X^2 \quad (\text{Property 1})$$

- Proof

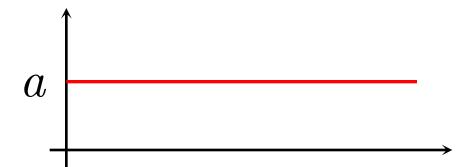
$$\begin{aligned} \text{Var}[X] &= E[(X - \mu_X)^2] = \sum_{x \in S_X} (x - \mu_X)^2 p_X(x) = \\ &= \sum_{x \in S_X} x^2 p_X(x) - 2\mu_X \underbrace{\sum_{x \in S_X} x p_X(x)}_{\mu_X} + \underbrace{\mu_X^2 \sum_{x \in S_X} p_X(x)}_1 = \\ &= E[X^2] - 2\mu_X^2 + \mu_X^2 = E[X^2] - \mu_X^2 \end{aligned}$$

- Variance of constant

$$\text{Var}[a] = 0 \quad (\text{Property 2})$$

- Proof

$$\text{Var}[a] = E[(a - E[a])^2] = E[0] = 0$$



# Recall of variance and its properties

---

- Invariance under the mean shift. Assume  $b$  is a constant.

$$\text{Var}[X + b] = \text{Var}[X] \quad (\text{Property 3})$$

- Proof

$$\begin{aligned}\text{Var}[X + b] &= E\left[\left((X + b) - E[X + b]\right)^2\right] = E\left[\left(X + b - E[X] - b\right)^2\right] \\ &= E\left[\left(X - E[X]\right)^2\right] = \text{Var}[X]\end{aligned}$$

- Scaling

$$\text{Var}[aX] = a^2 \text{Var}[X] \quad (\text{Property 4})$$

- Proof

$$\begin{aligned}\text{Var}[aX] &= E\left[\left(aX - E[aX]\right)^2\right] = E\left[\left(aX - aE[X]\right)^2\right] \\ &= E\left[a^2\left(X - E[X]\right)^2\right] = a^2 \text{Var}[X]\end{aligned}$$

# Recall of variance and its properties

---

- Example of covariance computation from joint distribution

## Example [edit]

Suppose that  $X$  and  $Y$  have the following joint probability mass function,<sup>[6]</sup> in which the six central cells give the discrete joint probabilities  $f(x, y)$  of the six hypothetical realizations  $(x, y) \in S = \{(1, 1), (1, 2), (1, 3), (2, 1), (2, 2), (2, 3)\}$  :

|   |           | <b>y</b> |     |     |
|---|-----------|----------|-----|-----|
|   | $f(x, y)$ | 1        | 2   | 3   |
| 1 |           | 1/4      | 1/4 | 0   |
| x | 2         | 0        | 1/4 | 1/4 |
|   | $f_Y(y)$  | 1/4      | 1/2 | 1/4 |

$X$  can take on two values (1 and 2) while  $Y$  can take on three (1, 2, and 3). Their means are  $\mu_X = 3/2$  and  $\mu_Y = 2$ . The population standard deviations of  $X$  and  $Y$  are  $\sigma_X = 1/2$  and  $\sigma_Y = \sqrt{2/3}$ . Then:

$$\begin{aligned}\text{cov}(X, Y) &= \sigma_{XY} = \sum_{(x,y) \in S} f(x, y)(x - \mu_X)(y - \mu_Y) \\ &= \left(\frac{1}{4}\right) \left(1 - \frac{3}{2}\right) (1 - 2) + \left(\frac{1}{4}\right) \left(1 - \frac{3}{2}\right) (2 - 2) \\ &\quad + (0) \left(1 - \frac{3}{2}\right) (3 - 2) + (0) \left(2 - \frac{3}{2}\right) (1 - 2) \\ &\quad + \left(\frac{1}{4}\right) \left(2 - \frac{3}{2}\right) (2 - 2) + \left(\frac{1}{4}\right) \left(2 - \frac{3}{2}\right) (3 - 2) \\ &= \frac{1}{4}.\end{aligned}$$

<https://en.wikipedia.org/wiki/Covariance>

# Appendix

---

- ▶ Recall B: matrix operations

# Recall of matrix operations: vector and matrix

## ■ Basic notations

### ▪ Vector

$$\mathbf{x} \in \mathbb{R}^n \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{matrix} & 1 & \textcircled{O} \\ & 2 & \textcircled{O} \\ & n & \textcircled{O} \end{matrix}$$

### ▪ Transposed vector

$$\mathbf{x}^T \in \mathbb{R}^{1 \times n} \triangleq [ - \quad \mathbf{x} \quad - ] = [ \ x_1 \quad x_2 \quad \cdots \quad x_n \ ]$$

### ▪ Matrix

$$\mathbf{A} \in \mathbb{R}^{n \times m} = \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ a_{12} & a_{22} & \cdots & a_{m2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1n} & a_{2n} & \cdots & a_{mn} \end{bmatrix}$$

# Recall of matrix operations: matrix forms

- Representation of matrix

- Column form

$$\mathbf{A} = \begin{bmatrix} & & & m \\ | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_m \\ | & | & & | \\ n & & & \end{bmatrix}$$

- Transposed column matrix

$$\mathbf{A}^T = \begin{bmatrix} & & n \\ - & \mathbf{a}_1^T & - \\ - & \mathbf{a}_2^T & - \\ \vdots & & \\ - & \mathbf{a}_m^T & - \\ m & & \end{bmatrix}$$

- Row form

$$\mathbf{A} = \begin{bmatrix} & & m \\ - & \vec{\mathbf{a}}_1^T & - \\ - & \vec{\mathbf{a}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{a}}_n^T & - \\ n & & \end{bmatrix}$$

- Transposed row matrix

$$\mathbf{A}^T = \begin{bmatrix} & & n \\ | & | & | \\ \vec{\mathbf{a}}_1 & \vec{\mathbf{a}}_2 & \cdots & \vec{\mathbf{a}}_n \\ | & | & & | \\ m & & & \end{bmatrix}$$

Note: we always assume that the vector is represented as a column.

# Recall of matrix operations: matrix multiplications

- **Vector-vector products**

- **Inner product**

$$\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} \in \mathbb{R} \triangleq \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} - & \mathbf{x} & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{y} \\ | \end{bmatrix} = \sum_{i=1}^n x_i y_i$$

- **Outer product**

$$\mathbf{x}\mathbf{y}^T \in \mathbb{R}^{m \times n} \triangleq \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \begin{bmatrix} y_1 & y_2 & \cdots & y_n \end{bmatrix} = \begin{bmatrix} x_1 y_1 & x_1 y_2 & \cdots & x_1 y_n \\ x_2 y_1 & x_2 y_2 & \cdots & x_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ x_m y_1 & x_m y_2 & \cdots & x_m y_n \end{bmatrix}$$

$$\mathbf{x}\mathbf{y}^T = \left[ \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} y_1 \quad \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} y_2 \quad \cdots \quad \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} y_n \right]$$

# Recall of matrix operations: matrix multiplications

- **Matrix-vector products**

$$\mathbf{y} = \mathbf{Ax} \in \mathbb{R}^m, \mathbf{x} \in \mathbb{R}^n, \mathbf{A} \in \mathbb{R}^{m \times n}$$

$$\mathbf{A} = \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix} = \begin{bmatrix} - & \vec{\mathbf{a}}_1^T & - \\ - & \vec{\mathbf{a}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{a}}_m^T & - \end{bmatrix}$$

- **Inner form**

$$\mathbf{y} = \mathbf{Ax} = \begin{bmatrix} - & \vec{\mathbf{a}}_1^T & - \\ - & \vec{\mathbf{a}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{a}}_m^T & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} \vec{\mathbf{a}}_1^T \mathbf{x} \\ \vec{\mathbf{a}}_2^T \mathbf{x} \\ \vdots \\ \vec{\mathbf{a}}_m^T \mathbf{x} \end{bmatrix}$$

- **Outer form**

$$\begin{aligned} \mathbf{y} = \mathbf{Ax} &= \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \\ &= \begin{bmatrix} | \\ \mathbf{a}_1 \\ | \end{bmatrix} x_1 + \cdots + \begin{bmatrix} | \\ \mathbf{a}_n \\ | \end{bmatrix} x_n \end{aligned}$$

# Recall of matrix operations: matrix multiplications

- Matrix-matrix products

$$\mathbf{AB}, \mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{m \times p}$$

$$\mathbf{A} = \begin{bmatrix} & & & & m \\ | & | & & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_m \\ | & | & & | \\ n \end{bmatrix} = \begin{bmatrix} & & & & m \\ - & \vec{\mathbf{a}}_1^T & - \\ - & \vec{\mathbf{a}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{a}}_n^T & - \\ n \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} & & & & p \\ | & | & & & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & | & & | \\ m \end{bmatrix} = \begin{bmatrix} & & & & p \\ - & \vec{\mathbf{b}}_1^T & - \\ - & \vec{\mathbf{b}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{b}}_m^T & - \\ m \end{bmatrix}$$

# Recall of matrix operations: matrix multiplications

- Matrix-matrix products

$$\mathbf{AB}, \mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{m \times p}$$

- Inner form

$$\begin{aligned}\mathbf{AB} &= \begin{bmatrix} - & \vec{\mathbf{a}}_1^T & - \\ - & \vec{\mathbf{a}}_2^T & - \\ \vdots & & \\ - & \vec{\mathbf{a}}_n^T & - \end{bmatrix} \begin{bmatrix} | & & & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & & & | \end{bmatrix} = \begin{bmatrix} \vec{\mathbf{a}}_1^T \mathbf{b}_1 & \vec{\mathbf{a}}_1^T \mathbf{b}_2 & \cdots & \vec{\mathbf{a}}_1^T \mathbf{b}_p \\ \vec{\mathbf{a}}_2^T \mathbf{b}_1 & \vec{\mathbf{a}}_2^T \mathbf{b}_2 & \cdots & \vec{\mathbf{a}}_2^T \mathbf{b}_p \\ \vdots & \vdots & \ddots & \vdots \\ \vec{\mathbf{a}}_n^T \mathbf{b}_1 & \vec{\mathbf{a}}_n^T \mathbf{b}_2 & \cdots & \vec{\mathbf{a}}_n^T \mathbf{b}_p \end{bmatrix} \\ &= \begin{bmatrix} | & & & | \\ \mathbf{Ab}_1 & \mathbf{Ab}_2 & \cdots & \mathbf{Ab}_p \\ | & & & | \end{bmatrix}\end{aligned}$$

# Recall of matrix operations: matrix multiplications

- Matrix-matrix products

$$\mathbf{AB}, \mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{m \times p}$$

- Outer form

$$\begin{aligned}\mathbf{AB} &= \left[ \begin{array}{cccc|c} & & & & & \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_m & & \\ & & & & & \\ n & & & & & \end{array} \right] \left[ \begin{array}{ccc|c} - & \vec{\mathbf{b}}_1^T & - & p \\ - & \vec{\mathbf{b}}_2^T & - & \\ \vdots & & & \\ - & \vec{\mathbf{b}}_m^T & - & \end{array} \right] \\ &= \underbrace{\left[ \begin{array}{c|cc} & & \\ \mathbf{a}_1 & - & \vec{\mathbf{b}}_1^T & - \\ & & & p \\ n & & & \end{array} \right]}_{n \times p} + \underbrace{\left[ \begin{array}{c|cc} & & \\ \mathbf{a}_2 & - & \vec{\mathbf{b}}_2^T & - \\ & & & p \\ n & & & \end{array} \right]}_{n \times p} + \cdots + \underbrace{\left[ \begin{array}{c|cc} & & \\ \mathbf{a}_m & - & \vec{\mathbf{b}}_m^T & - \\ & & & p \\ n & & & \end{array} \right]}_{n \times p}\end{aligned}$$

# Applications: sample covariance matrix

- Sample covariance matrix (assuming zero-mean)

$$\mathbf{C} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{n \times n}, \mathbf{X} \in \mathbb{R}^{n \times N}$$

$$\mathbf{C} = \frac{1}{N} \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ | & | & & | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ \vdots & & \\ - & \mathbf{x}_N^T & - \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ | & | & & | \end{bmatrix}$$

Assume zero mean vectors

- Outer form

$$\mathbf{C} = \frac{1}{N} \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ | & | & & | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ \vdots & & \\ - & \mathbf{x}_N^T & - \end{bmatrix} = \frac{1}{N} [\mathbf{x}_1 \mathbf{x}_1^T + \mathbf{x}_2 \mathbf{x}_2^T + \cdots + \mathbf{x}_N \mathbf{x}_N^T]$$

$$= \underbrace{\begin{bmatrix} | \\ \mathbf{x}_1 \\ | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_1^T & - \end{bmatrix}}_{n \times n} + \underbrace{\begin{bmatrix} | \\ \mathbf{x}_2 \\ | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_2^T & - \end{bmatrix}}_{n \times n} + \cdots + \underbrace{\begin{bmatrix} | \\ \mathbf{x}_N \\ | \end{bmatrix} \begin{bmatrix} - & \mathbf{x}_N^T & - \end{bmatrix}}_{n \times n}$$

# Applications: Gram matrix

- **Gram matrix**

$$\mathbf{G} = \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{N \times N}, \mathbf{X} \in \mathbb{R}^{n \times N}$$

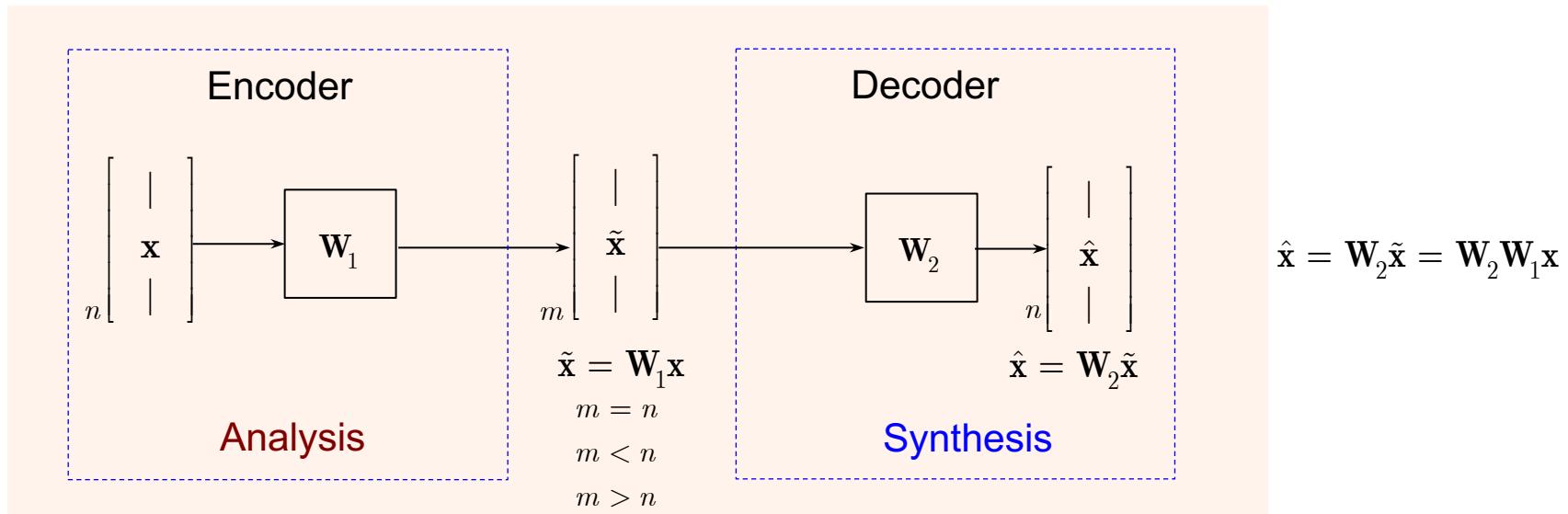
$$\mathbf{X} = \begin{bmatrix} & & & \\ | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ | & | & & | \\ n & & & \end{bmatrix}$$

- **Inner form**

$$\mathbf{C} = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ \vdots & & \\ N & \mathbf{x}_N^T & - \end{bmatrix} \begin{bmatrix} & & & \\ | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_N \\ | & | & & | \\ & & & N \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \mathbf{x}_1 & \mathbf{x}_1^T \mathbf{x}_2 & \cdots & \mathbf{x}_1^T \mathbf{x}_N \\ \mathbf{x}_2^T \mathbf{x}_1 & \mathbf{x}_2^T \mathbf{x}_2 & \cdots & \mathbf{x}_2^T \mathbf{x}_N \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_N^T \mathbf{x}_1 & \mathbf{x}_N^T \mathbf{x}_2 & \cdots & \mathbf{x}_N^T \mathbf{x}_N \end{bmatrix}$$

# Applications: encoder-decoder notation

- Encoder-decoder notations



- Reconstruction (if no other constraints on properties of  $\tilde{\mathbf{x}}$ ) for given  $\mathbf{W}_1$  and  $\mathbf{W}_2$

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}\|_2^2 = \|(\mathbf{I} - \mathbf{W}_2 \mathbf{W}_1) \mathbf{x}\|_2^2 \Rightarrow \text{For perfect reconstruction}$$
$$(\mathbf{I} - \mathbf{W}_2 \mathbf{W}_1) = \mathbf{0} \rightarrow \mathbf{I} = \mathbf{W}_2 \mathbf{W}_1$$

# Applications: encoder-decoder notation

- Encoder-decoder notations

- Reconstruction

$$\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2 = \|\mathbf{x} - \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}\|_2^2 = \|(\mathbf{I} - \mathbf{W}_2 \mathbf{W}_1) \mathbf{x}\|_2^2 \Rightarrow \mathbf{W}_2 \mathbf{W}_1 = \mathbf{I}$$

- Special case  $m = n$ ,  $\mathbf{W}_2 = \mathbf{A} \in \mathbb{R}^{n \times n}$  = 
$$\begin{bmatrix} | & | & & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_n \\ | & | & & | \end{bmatrix}$$
 - orthonormal  
Orthonormal matrices  $\mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T = \mathbf{I}$

$$\mathbf{W}_1 = \mathbf{A}^T \in \mathbb{R}^{n \times n} = \begin{bmatrix} - & \mathbf{a}_1^T & - \\ - & \mathbf{a}_2^T & - \\ \vdots & & \\ - & \mathbf{a}_n^T & - \end{bmatrix} \quad \begin{array}{l} \mathbf{a}_i^T \mathbf{a}_i = 1, \\ \mathbf{a}_i^T \mathbf{a}_j = 0 \end{array}$$

$$\hat{\mathbf{x}} = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} = \underbrace{\mathbf{A} \mathbf{A}^T}_{\mathbf{I}} \mathbf{x} = \mathbf{x}$$

- Otherwise, when  $m \neq n$ , a special care should be taken about the selection of encoder-decoder pairs. We will see more details in pseudo-inverse part.

# Appendix

---

- ▶ Recall C: eigen value decomposition and SVD

# Eigen (value) decomposition

- Summary (if you are experienced)

- $\mathbf{X}$  is a **square real** matrix

- Eigenvectors

$$\mathbf{U} \in \underbrace{\mathbb{C}^{n \times n}}_{\text{generally complex}} = \begin{bmatrix} | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_n \\ | & | & & | \end{bmatrix} \quad \mathbf{U}\mathbf{U}^{-1} = \mathbf{I}$$

- Eigenvalues  $\sigma_1, \dots, \sigma_n$ ,  $\underbrace{\sigma_i \in \mathbb{C}}_{\text{generally complex}}$

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n) \text{ or } \Sigma_{ii} = \sigma_i, \forall i$$

- Eigen decomposition

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{U}^{-1}$$

# Eigen (value) decomposition

- **Summary (if you are experienced)**

- $X$  is a (square) symmetric real matrix

- Eigenvectors are orthonormal

$$U^{-1} = U^T$$

- Eigenvalues  $\sigma_1, \dots, \sigma_n$  are real  $\sigma_i \in \mathbb{R}$

- Eigen decomposition of symmetric real matrix

$$X = U\Sigma U^T$$

$$X \in \mathbb{R}^{n \times n} = \begin{bmatrix} a & \times & & \cdots & \circ \\ \times & b & & & \\ & & c & & \\ \vdots & & & \ddots & \\ \circ & & & & z \end{bmatrix}$$

## Matlab

```
[U, D] = eig(A)
```

# Eigen (value) decomposition (you can skip it)

- **Friendly introduction**

- X is any square real matrix

- **Definition:** the eigenvector  $\mathbf{u}$  is a vector that satisfies:

$$\mathbf{X}\mathbf{u} = \sigma\mathbf{u}, \mathbf{u} \neq \mathbf{0} \Rightarrow (\mathbf{X} - \sigma\mathbf{I})\mathbf{u} = \mathbf{0}$$

Example

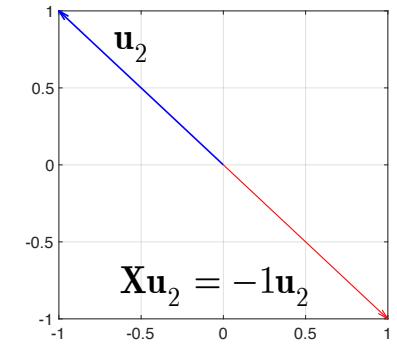
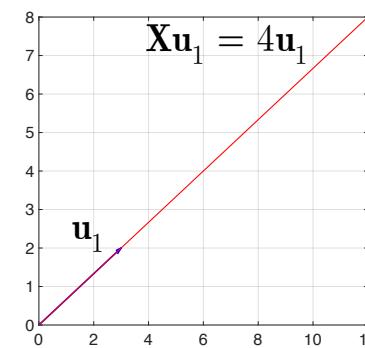
$\mathbf{X} = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix}$  has the eigenvectors

$$\mathbf{u}_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}, \text{ with } \sigma_1 = 4$$

$$\mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \text{ with } \sigma_2 = -1$$

$$\mathbf{X}\mathbf{u}_1 = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \underbrace{4}_{\sigma_1} \begin{bmatrix} 3 \\ 2 \end{bmatrix} = \begin{bmatrix} 12 \\ 8 \end{bmatrix}$$

$$\mathbf{X}\mathbf{u}_2 = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \underbrace{-1}_{\sigma_2} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$



# Eigen (value) decomposition (you can skip it)

---

- **Friendly introduction**

- In most applications, the eigenvectors are normalized

$$\mathbf{u}_i^T \mathbf{u}_i = 1$$

## Example

$$\mathbf{u}_1 = \begin{bmatrix} 0.83 \\ 0.55 \end{bmatrix}, \text{ with } \sigma_1=4$$

$$\mathbf{u}_2 = \begin{bmatrix} -0.71 \\ 0.71 \end{bmatrix}, \text{ with } \sigma_2=-1$$

## Matlab

```
>> X=[2 3; 2 1]
X=
      2      3
      2      1
>> [U,D] = eig(X)
U =
      0.8321    -0.7071
      0.5547     0.7071
D =
      4      0
      0     -1
```

# Eigen (value) decomposition (you can skip it)

- **Friendly introduction**

- **Definition:** matrix form

$$\mathbf{X}\mathbf{u}_i = \sigma_i \mathbf{u}_i, \mathbf{u}_i \neq \mathbf{0}, \forall i \Rightarrow \mathbf{X}\mathbf{U} = \mathbf{U}\Sigma$$

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{U}^{-1}$$

Example: to check it out that indeed  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{U}^{-1}$

## Matlab

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{U}^{-1}$$

$$= \begin{bmatrix} 3 & -1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 0.2 & 0.2 \\ -0.4 & 0.6 \end{bmatrix}$$
$$= \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix}$$

```
>> res=U*D*inv(U)  
res =  
2.0000 3.0000  
2.0000 1.0000
```

# Eigen (value) decomposition (you can skip it)

## ▪ Friendly introduction

- **Definition 1:** positive (semi-)definite matrices
  - A matrix is said to be positive semi-definite when it can be obtained as the product of a matrix by its transpose

$$\mathbf{C} = \mathbf{XX}^T$$

- This implies that the positive semi-definite matrix is **always symmetric**.

- **Definition 2 (alternative):** positive (semi-)definite matrices
  - We can investigate the sign of eigenvalues of symmetric  $\mathbf{X}$
  - Consider a decomposition  $\mathbf{a}^T \mathbf{X} \mathbf{a}$  for any non-zero vector  $\mathbf{a}$

$$\mathbf{a}^T \mathbf{X} \mathbf{a} = \underbrace{\mathbf{a}^T \mathbf{U}}_{\mathbf{y}^T = \mathbf{a}^T \mathbf{U}} \Sigma \underbrace{\mathbf{U}^T \mathbf{a}}_{\mathbf{y} = \mathbf{U}^T \mathbf{a}} = \sum_{i=1}^n \sigma_i \underbrace{y_i^2}_{\geq 0}$$

Since  $y_i^2$  is always positive, the sign of the above quadratic form is fully determined by the sign of eigenvalues.

# Eigen (value) decomposition (you can skip it)

---

- **Friendly introduction**

- **Definition 2:** positive (semi-)definite matrices
  - If all  $\sigma_i \geq 0$ , then the matrix  $X$  is said to be **positive semi-definite**
  - If all  $\sigma_i > 0$ , then the matrix  $X$  is said to be **positive definite**
  - If all  $\sigma_i \leq 0$ , then the matrix  $X$  is said to be **negative semi-definite**
  - If all  $\sigma_i < 0$ , then the matrix  $X$  is said to be **negative definite**

# Eigen (value) decomposition (you can skip it)

---

- **Friendly introduction**
  - Properties of positive (semi-)definite matrices
    - Its eigenvalues are always positive or null, i.e.,  $\sigma_i \geq 0$
    - The eigenvectors are pairwise orthonormal, if their eigenvalues are different
    - The eigenvectors are real
    - Since all eigenvectors are orthonormal, i.e.,  $\mathbf{U}^{-1} = \mathbf{U}^T$
  - Rank of matrix:
    - the number of non-zero eigenvalues of the matrix and  $r \leq \min \{ n, m \}$
    - If rank is equal to the matrix dimension, it is called a **full rank** matrix
    - Otherwise, it is called **rank-deficient**, **singular** or **multicollinear** matrix
  - To have a full rank, all vectors should be linearly independent, i.e., none of them can be created as a linear combination of remaining ones

# Singular value decomposition (SVD)

- **Fundamentals of SVD (can be skipped, if familiar)**

- For any  $\mathbf{X} \in \mathbb{R}^{n \times m}$ , there exists a factorization  $\underbrace{\mathbf{X}}_{n \times m} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\mathbf{S}}_{n \times m} \underbrace{\mathbf{V}^T}_{m \times m}$

$$n \begin{bmatrix} & & & & m \\ | & & & & | \\ \mathbf{x}_1 & \cdots & \mathbf{x}_m \\ | & & | \\ & & & & n \end{bmatrix} = \underbrace{n \begin{bmatrix} & & & & \\ | & & | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_r & \mathbf{u}_{r+1} & \cdots & \mathbf{u}_n \\ | & & | & | & & | \\ & & & & & n \end{bmatrix}}_{\mathbf{U}} \underbrace{\begin{bmatrix} \lambda_1 & & & & & & m \\ & \ddots & & & & & | \\ & & \lambda_r & & & & | \\ & & & 0 & & & | \\ & & & & \ddots & & | \\ & & & & & 0 & & | \\ & & & & & & & n \end{bmatrix}}_{\mathbf{S}} \underbrace{\begin{bmatrix} & & & & & & m \\ - & \mathbf{v}_1^T & - & & & & | \\ - & \mathbf{v}_r^T & - & & & & | \\ - & \mathbf{v}_{r+1}^T & - & & & & | \\ - & \mathbf{v}_m^T & - & & & & | \\ & & & & & & - \end{bmatrix}}_{\mathbf{V}^T}$$

- $\mathbf{U}$  is an  $n \times n$  orthonormal matrix consisting of  $\mathbf{U} = \text{evec}(\mathbf{XX}^T)$
- $\mathbf{V}$  is an  $m \times m$  orthonormal matrix consisting of  $\mathbf{V} = \text{evec}(\mathbf{X}^T \mathbf{X})$
- $\mathbf{S}$  is an  $n \times m$  diagonal matrix containing  $r \leq \min\{n, m\}$  singular values  $\lambda_i \geq 0$  on the main diagonal and  $\mathbf{S}^2 = \text{eval}(\mathbf{X}^T \mathbf{X}) = \text{eval}(\mathbf{XX}^T)$

# Singular value decomposition (SVD)

- Properties of SVD:
  - $\underbrace{\mathbf{XX}^T}_{n \times n}$  and  $\underbrace{\mathbf{X}^T \mathbf{X}}_{m \times m}$  are **square symmetric** matrices that can be decomped as:

EVD: 
$$\begin{aligned}\mathbf{XX}^T &= \mathbf{U}\Sigma\mathbf{U}^T & \Sigma &= \text{eval}(\mathbf{XX}^T) \\ \mathbf{X}^T \mathbf{X} &= \mathbf{V}\Sigma'\mathbf{V}^T & \text{where} & \Sigma' = \text{eval}(\mathbf{X}^T \mathbf{X})\end{aligned}$$

$\Sigma$  and  $\Sigma'$  have the same non-zero diagonal entries  
but in the different order

Validate:  $\mathbf{XX}^T = \mathbf{USV}^T (\mathbf{USV}^T)^T = \mathbf{US} \underbrace{\mathbf{V}^T \mathbf{V}}_{\mathbf{I}} \mathbf{S}^T \mathbf{U}^T = \mathbf{U} \underbrace{\mathbf{SS}^T}_{\mathbf{S}^2 = \Sigma, \lambda_i^2 = \sigma_i} \mathbf{U}^T = \mathbf{U}\Sigma\mathbf{U}^T$

# Singular value decomposition (SVD)

---

- **Demo example**

X =

$$\begin{matrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{matrix}$$

>> [U, S, V] = svd(A)

U =

$$\begin{matrix} -0.7071 & -0.7071 \\ -0.7071 & 0.7071 \end{matrix}$$

S =

$$\lambda_i \begin{matrix} 3.4641 & 0 & 0 \\ 0 & 3.1623 & 0 \end{matrix}$$

V =

$$\begin{matrix} -0.4082 & -0.8944 & -0.1826 \\ -0.8165 & 0.4472 & -0.3651 \\ -0.4082 & 0.0000 & 0.9129 \end{matrix}$$

# Singular value decomposition (SVD)

## ▪ Fundamentals of SVD

### ▪ Demo example

$X =$

$$\begin{matrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{matrix}$$

```
>> [U, D] = eig(X*X')
```

$U =$

$$\begin{matrix} -0.7071 & 0.7071 \\ 0.7071 & 0.7071 \end{matrix}$$

$D =$

$$\begin{bmatrix} 10.0000 & 0 \\ 0 & 12.0000 \end{bmatrix}$$

1. Order

2. Squared

$$S^2 = \Sigma, \lambda_i^2 = \sigma_i$$

Symmetric

```
>> X*X'
```

$ans =$

$$\begin{matrix} 11 & 1 \\ 1 & 11 \end{matrix}$$

```
>> X'*X
```

$ans =$

$$\begin{matrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{matrix}$$

```
>> [V, D] = eig(X'*X)
```

$V =$

$$\begin{matrix} 0.1826 & 0.8944 & 0.4082 \\ 0.3651 & -0.4472 & 0.8165 \\ -0.9129 & 0 & 0.4082 \end{matrix}$$

$D =$

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 10 & 0 \\ 0 & 0 & 12 \end{bmatrix}$$

# Singular value decomposition (SVD)

- Economy sized SVD or thin SVD**

- Avoids computing unnecessary components

$$\underbrace{\mathbf{X}}_{n \times m} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\mathbf{S}}_{n \times m} \underbrace{\mathbf{V}^T}_{m \times m} \longrightarrow \underbrace{\mathbf{X}}_{n \times m} = \underbrace{\tilde{\mathbf{U}}}_{n \times r} \underbrace{\tilde{\mathbf{S}}}_{r \times r} \underbrace{\tilde{\mathbf{V}}^T}_{r \times m}$$

$$n \begin{bmatrix} | & & m \\ \mathbf{x}_1 & \cdots & \mathbf{x}_m \\ | & & | \end{bmatrix} = \underbrace{\begin{bmatrix} | & & | \\ \mathbf{u}_1 & \cdots & \mathbf{u}_r \\ | & & | \end{bmatrix}}_{\mathbf{U}} \begin{bmatrix} | & & | \\ \mathbf{u}_{r+1} & \cdots & \mathbf{u}_n \\ | & & | \end{bmatrix} \underbrace{\begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_r & \\ & & & 0 \end{bmatrix}}_S \underbrace{\begin{bmatrix} m \\ - & \mathbf{v}_1^T & - \\ - & \mathbf{v}_r^T & - \\ - & \mathbf{v}_{r+1}^T & - \\ - & \mathbf{v}_m^T & - \end{bmatrix}}_{\mathbf{V}^T}$$

Since there is at most  $r \leq \min\{n, m\}$  non-zero singular values.

- Truncated SVD**

- Shrink with  $k < r$ :

$$\underbrace{\mathbf{X}}_{n \times m} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\mathbf{S}}_{n \times m} \underbrace{\mathbf{V}^T}_{m \times m} \longrightarrow \underbrace{\mathbf{X}}_{n \times m} = \underbrace{\tilde{\mathbf{U}}}_{n \times k} \underbrace{\tilde{\mathbf{S}}}_{k \times k} \underbrace{\tilde{\mathbf{V}}^T}_{k \times m}$$

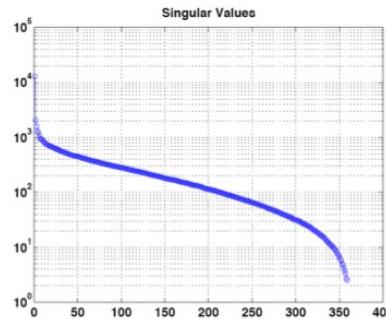
# Singular value decomposition (SVD)

- **Spectral view on SVD**

- Spectrum of  $\mathbf{X}$   $\{\lambda_i\}_{i=1}^r$  decays very quickly and these singular values are ordered

$$\underbrace{\mathbf{X}}_{n \times m} = \underbrace{\tilde{\mathbf{U}}}_{n \times k} \underbrace{\tilde{\mathbf{S}}}_{k \times k} \underbrace{\tilde{\mathbf{V}}^T}_{k \times m} = \sum_{i=1}^k \lambda_i \mathbf{u}_i \mathbf{v}_i^T$$

$$\mathbf{X} = [\lambda_1] \underbrace{\begin{bmatrix} | \\ \mathbf{u}_1 \\ | \end{bmatrix} \left[ \begin{array}{ccc} - & \mathbf{v}_1^T & - \end{array} \right]}_{n \times m} + [\lambda_2] \underbrace{\begin{bmatrix} | \\ \mathbf{u}_2 \\ | \end{bmatrix} \left[ \begin{array}{ccc} - & \mathbf{v}_2^T & - \end{array} \right]}_{n \times m} + \dots + [\lambda_k] \underbrace{\begin{bmatrix} | \\ \mathbf{u}_k \\ | \end{bmatrix} \left[ \begin{array}{ccc} - & \mathbf{v}_k^T & - \end{array} \right]}_{n \times m}$$



# Singular value decomposition (SVD)

- Case  $m > n$ : economy SVD

$$n \begin{array}{|c|} \hline m \\ \hline \end{array} X = n \begin{array}{|c|} \hline n \\ \hline \end{array} U \begin{array}{|c|} \hline n \\ \hline \end{array} \begin{array}{|c|} \hline m \\ \hline \end{array} S \begin{array}{|c|} \hline m \\ \hline \end{array} V^T$$
$$n \begin{array}{|c|} \hline m \\ \hline \end{array} X \approx n \begin{array}{|c|} \hline n \\ \hline \end{array} \tilde{U} \begin{array}{|c|} \hline n \\ \hline \end{array} \begin{array}{|c|} \hline m \\ \hline \end{array} \tilde{V}^T$$

$$U = \tilde{U} \quad \tilde{U}^T \tilde{U} = \tilde{U} \tilde{U}^T = I_n$$

$$\tilde{V}^T \tilde{V} = I \neq \tilde{V} \tilde{V}^T$$

# Singular value decomposition (SVD)

- Case  $n > m$ : economy SVD

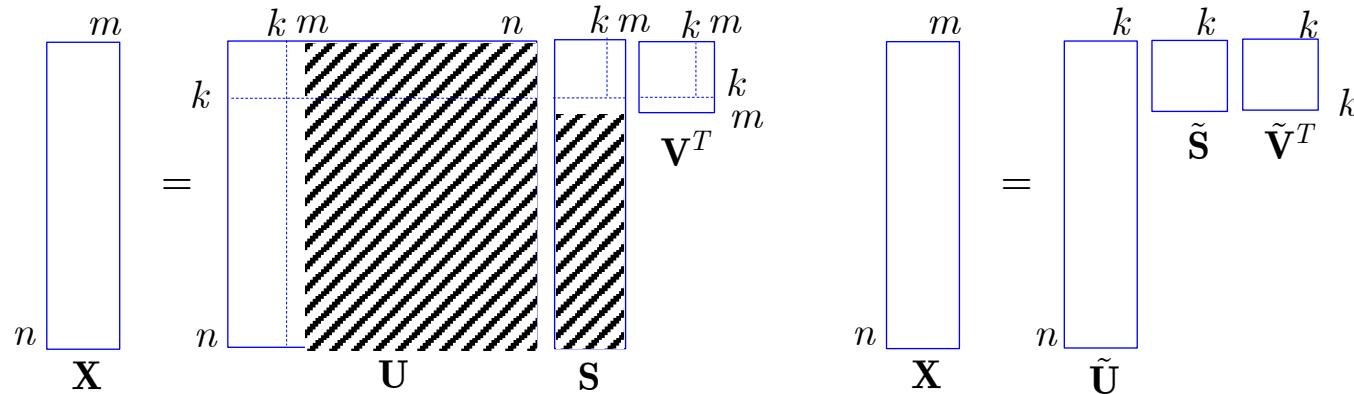
$$\begin{matrix} m \\ n \end{matrix} \begin{matrix} X \end{matrix} = \begin{matrix} m \\ n \end{matrix} \begin{matrix} U \end{matrix} \begin{matrix} n \\ m \\ m \end{matrix} \begin{matrix} S \end{matrix} \begin{matrix} m \\ m \end{matrix} \begin{matrix} V^T \end{matrix}$$
$$\begin{matrix} m \\ n \end{matrix} \begin{matrix} X \end{matrix} = \begin{matrix} m \\ n \end{matrix} \begin{matrix} \tilde{U} \end{matrix} \begin{matrix} m \\ m \\ m \end{matrix} \begin{matrix} \tilde{S} \end{matrix} \begin{matrix} \tilde{V}^T \end{matrix}$$

$$V = \tilde{V} \quad \tilde{V}^T \tilde{V} = \tilde{V} \tilde{V}^T = I_m$$

$$\tilde{U}^T \tilde{U} = I \neq \tilde{U} \tilde{U}^T$$

# Singular value decomposition (SVD)

- Case  $k < \min\{m, n\}$ : truncated SVD



$$\tilde{\mathbf{V}}^T \tilde{\mathbf{V}} = \mathbf{I} \neq \tilde{\mathbf{V}} \tilde{\mathbf{V}}^T$$

$$\tilde{\mathbf{U}}^T \tilde{\mathbf{U}} = \mathbf{I} \neq \tilde{\mathbf{U}} \tilde{\mathbf{U}}^T$$

# Notations

- **Characteristic function (CF)** of pdf is defined as:

$$\Phi(t) = \int_{-\infty}^{\infty} f(x)e^{jtx}dx = \mathbb{E}_{f(x)}[e^{jtx}]$$

where  $j = \sqrt{-1}$

- The pdf can be recovered as:

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \Phi(t)e^{-jtx}dt$$

**Moment-generating** function of pdf:

$$M(t) = \mathbb{E}_{f(x)}[e^{tX}]$$

# Appendix

---

- ▶ Recall D: moments of pdf and cdf

# Notations

---

- Taylor's series expansion:  $e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots = \sum_{k=0}^{\infty} \frac{x^k}{k!}$

$$\exp(tX) = \sum_{n=0}^{\infty} \frac{t^n \cdot X^n}{n!}$$

$$\exp(jtX) = \sum_{n=0}^{\infty} \frac{(jt)^n \cdot X^n}{n!}$$

$$\Phi(t) = \mathbb{E}_{f(x)}[\exp(jtX)] = 1 + \frac{jt\mathbb{E}_{f(x)}[X]}{1} - \frac{t^2\mathbb{E}_{f(x)}[X^2]}{2!} + \dots + \frac{(jt)^n\mathbb{E}_{f(x)}[X^n]}{n!}$$

$$M(t) = \mathbb{E}_{f(x)}[\exp(tX)] = 1 + \frac{t\mathbb{E}_{f(x)}[X]}{1} + \frac{t^2\mathbb{E}_{f(x)}[X^2]}{2!} + \dots + \frac{t^n\mathbb{E}_{f(x)}[X^n]}{n!}$$

The first four moments define the mean, variance, skewness, and kurtosis of the PDF  $f(x)$ , respectively.

# Notations

---

- Practical computation of moments from i.i.d. sequence  $\mathbf{x} = (x_1, \dots, x_N)$  drawn from  $f(x)$

- The nth **empirical pdf moment estimation**

$$\hat{m}_n = \frac{1}{N} \sum_{i=1}^N x_i^n, \quad n \geq 1$$

- Which is an unbiased estimation of **true pdf moments**

$$m_n = \mathbb{E}[X^n] = \int_{-\infty}^{\infty} f(x) x^n dx$$

- **Moments of CF:**

$$\mathbf{M}_n = \int_{-\infty}^{\infty} \Phi(t) t^n dt$$

$$\text{Link to pdf: } \mathbf{M}_n = j^n 2\pi \left. \frac{d^n}{dx^n} f(x) \right|_{x=0}$$

# Notations

---

- Practical computation of moments from i.i.d. sequence  $\mathbf{x} = (x_1, \dots, x_N)$  drawn from  $f(x)$ 
  - **Empirical moments of CDF:**

- Estimate empirical M-bin histogram  $\{h(m)\}_{m=0}^{M-1}$  of pdf  $f(x)$
- Estimate K-point CDF with  $K = 2^{\lceil \log_2 M \rceil}$  as

$$\Phi(k) = \sum_{m=0}^{M-1} h(m) \exp \left\{ \frac{j2\pi mk}{K} \right\}$$

Note that Fourier works in two ways:

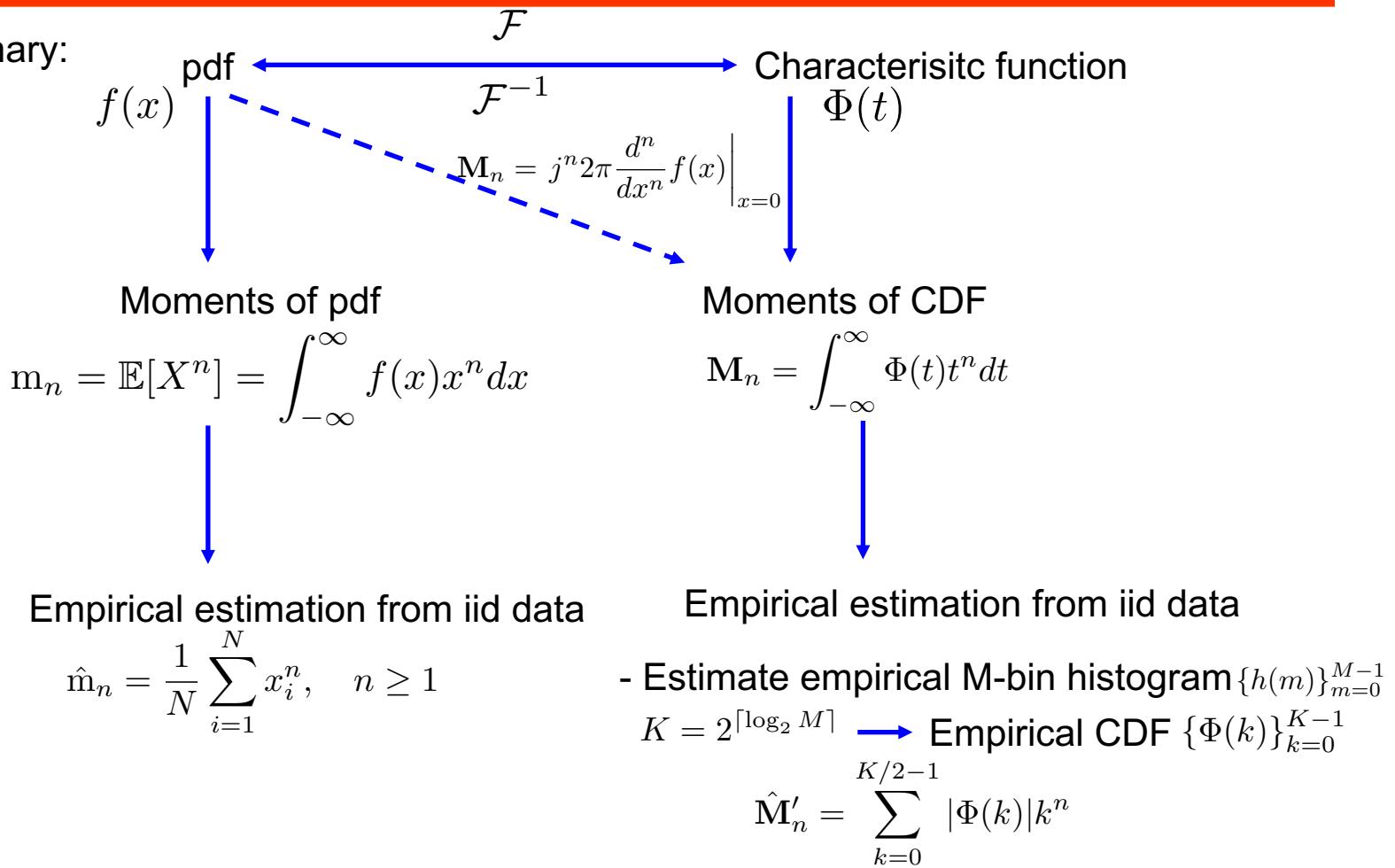
$$h(m) = \frac{1}{K} \sum_{k=0}^{K-1} \Phi(k) \exp \left\{ -\frac{j2\pi mk}{K} \right\}$$
$$0 \leq m \leq M-1$$

- Compute empirical CDF moments as:

$$\hat{\mathbf{M}}'_n = \sum_{k=0}^{K/2-1} |\Phi(k)| k^n \quad \hat{\mathbf{M}}_n = \sum_{k=0}^{K-1} \Phi(k) \sin^n \left( \frac{\pi k}{K} \right)$$

# Notations

- Summary:



# Notations

Applications:

- Outlier detection
- Fake detection
- Secret message detection in images (steganalysis)

